

# Applied Machine Learning

Regularization

**Siamak Ravanbakhsh**

COMP 551 (winter 2020)

# Learning objectives

Basic idea of

- overfitting and underfitting
- Regularization (L1 & L2)
- MLE vs MAP estimation
- bias and variance trade off
- evaluation metrics & cross validation

# Previously...

Linear regression and logistic regression  
is linear too simple? what if it's not a good fit?  
how to increase the models expressiveness?

- create new nonlinear features
- is there a downside?

# Recall: nonlinear basis functions

replace original features in  $f_w(x) = \sum_d w_d x_d$

# Recall: nonlinear basis functions

replace original features in  $f_w(x) = \sum_d w_d x_d$

with nonlinear bases  $f_w(x) = \sum_d w_d \phi_d(x)$

# Recall: nonlinear basis functions

replace original features in  $f_w(x) = \sum_d w_d x_d$

with nonlinear bases  $f_w(x) = \sum_d w_d \phi_d(x)$

linear least squares solution  $w^* = (\Phi^\top \Phi)^{-1} \Phi^\top y$

replacing  $X$  with  $\Phi$

$$\Phi = \begin{bmatrix} \phi_1(x^{(1)}), & \phi_2(x^{(1)}), & \cdots, & \phi_D(x^{(1)}) \\ \phi_1(x^{(2)}), & \phi_2(x^{(2)}), & \cdots, & \phi_D(x^{(2)}) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_1(x^{(N)}), & \phi_2(x^{(N)}), & \cdots, & \phi_D(x^{(N)}) \end{bmatrix}$$

# Recall: nonlinear basis functions

replace original features in  $f_w(x) = \sum_d w_d x_d$

with nonlinear bases  $f_w(x) = \sum_d w_d \phi_d(x)$

linear least squares solution  $w^* = (\Phi^\top \Phi)^{-1} \Phi^\top y$

replacing  $X$  with  $\Phi$

a (nonlinear) feature

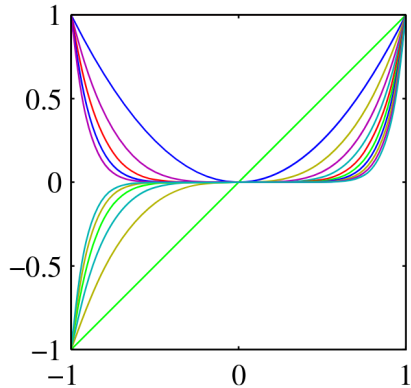
$$\Phi = \begin{bmatrix} \phi_1(x^{(1)}), & \phi_2(x^{(1)}), & \cdots, & \phi_D(x^{(1)}) \\ \phi_1(x^{(2)}), & \phi_2(x^{(2)}), & \cdots, & \phi_D(x^{(2)}) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_1(x^{(N)}), & \phi_2(x^{(N)}), & \cdots, & \phi_D(x^{(N)}) \end{bmatrix}$$

one instance

# Recall: nonlinear basis functions

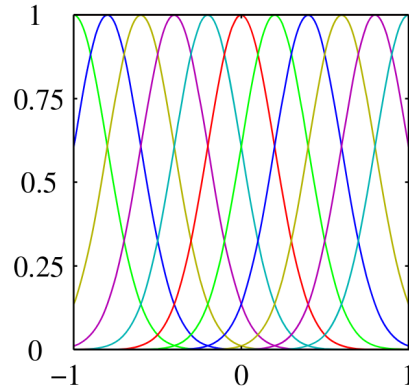
examples

original input is scalar  $x \in \mathbb{R}$



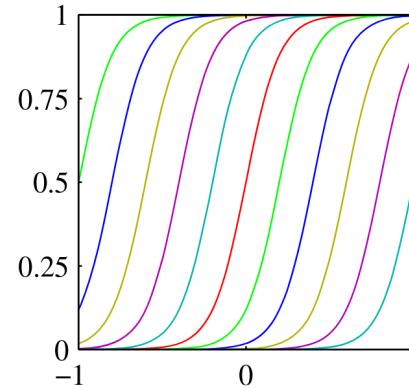
polynomial bases

$$\phi_k(x) = x^k$$



Gaussian bases

$$\phi_k(x) = e^{-\frac{(x-\mu_k)^2}{s^2}}$$

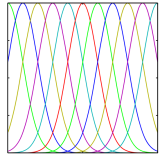


Sigmoid bases

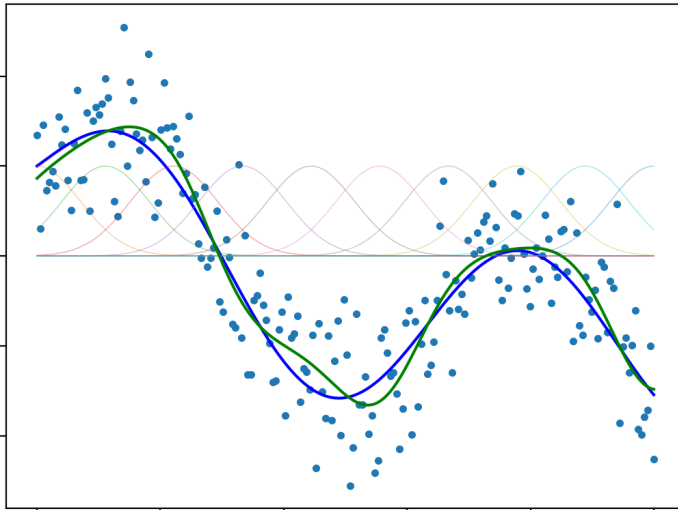
$$\phi_k(x) = \frac{1}{1+e^{-\frac{x-\mu_k}{s}}}$$



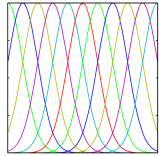
# Example: Gaussian bases



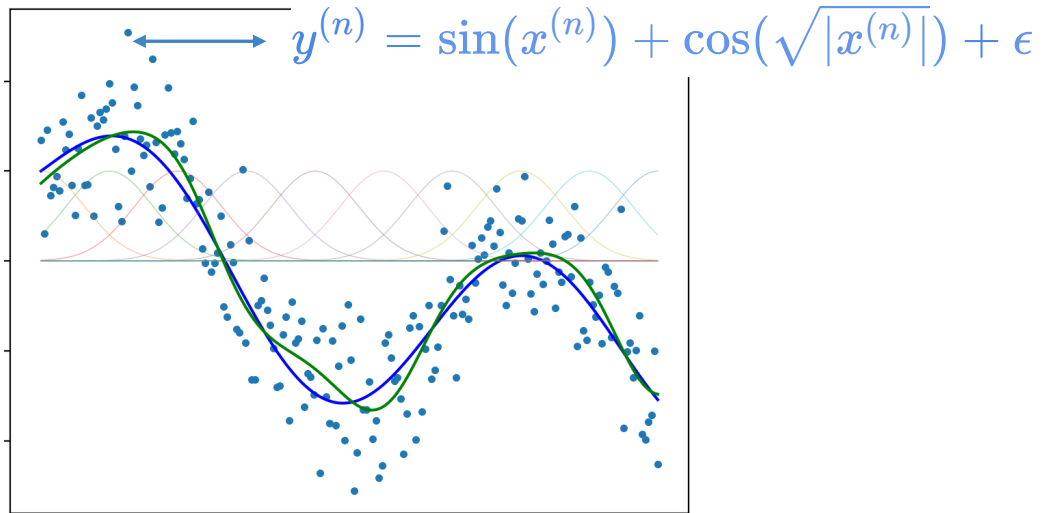
$$\phi_k(x) = e^{-\frac{(x-\mu_k)^2}{s^2}}$$



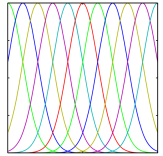
# Example: Gaussian bases



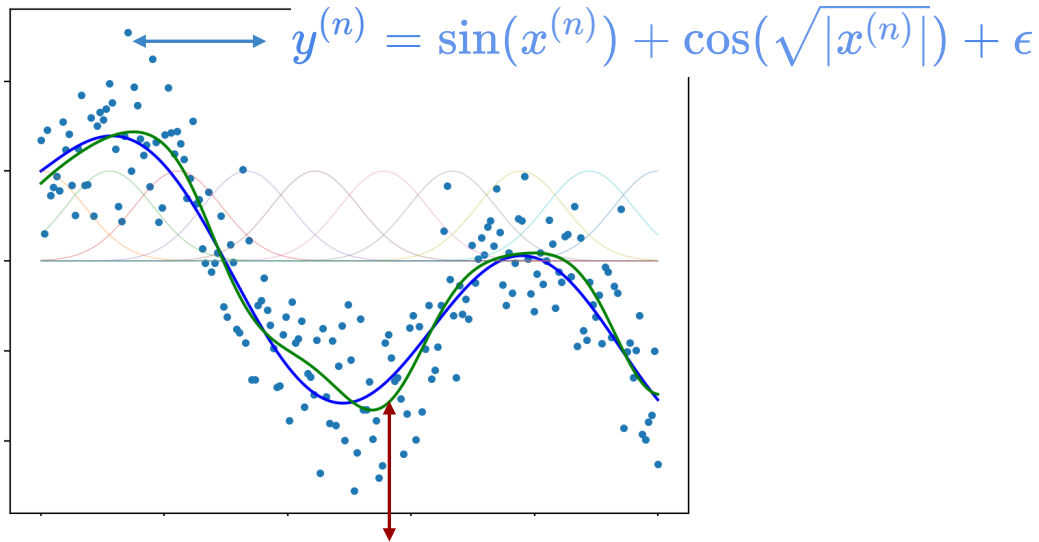
$$\phi_k(x) = e^{-\frac{(x-\mu_k)^2}{s^2}}$$



# Example: Gaussian bases

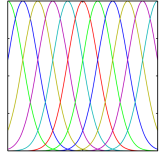


$$\phi_k(x) = e^{-\frac{(x-\mu_k)^2}{s^2}}$$

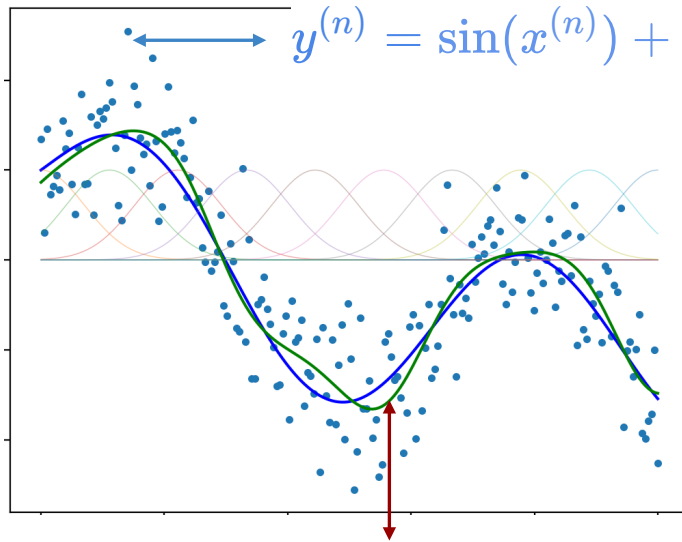


our fit to data using 10 Gaussian bases

# Example: Gaussian bases



$$\phi_k(x) = e^{-\frac{(x-\mu_k)^2}{s^2}}$$



our fit to data using 10 Gaussian bases

prediction for a new instance

$$f(x') = \phi(x')^\top (\Phi^\top \Phi)^{-1} \Phi^\top y$$

new instance

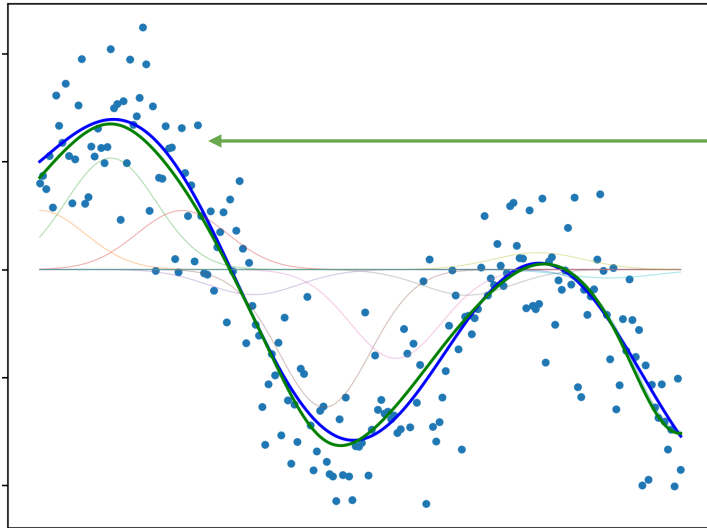
features evaluated for the new point

$w$  found using LLS

# Example: Gaussian bases



$$\phi_k(x) = e^{-\frac{(x-\mu_k)^2}{s^2}}$$

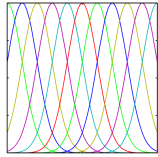


our fit to data using **10 Gaussian bases**

why not more?

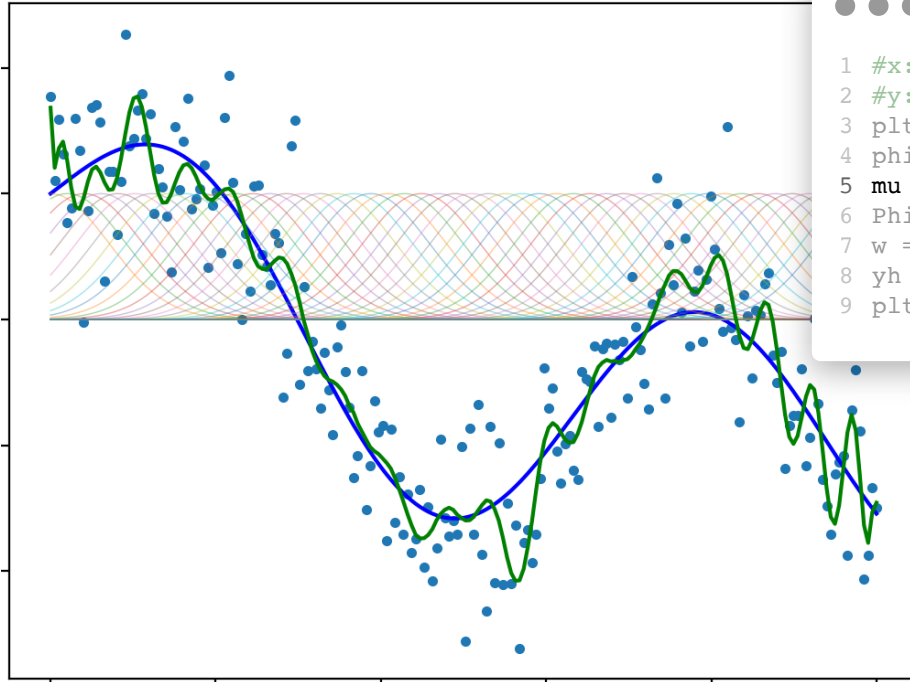
```
1 #x: N
2 #y: N
3 plt.plot(x, y, 'b.')
4 phi = lambda x,mu: np.exp(-(x-mu)**2)
5 mu = np.linspace(0,10,10) #10 Gaussian bases
6 Phi = phi(x[:,None], mu[None,:]) #N x 10
7 w = np.linalg.lstsq(Phi, y)[0]
8 yh = np.dot(Phi,w)
9 plt.plot(x, yh, 'g-')
```

# Example: Gaussian bases



$$\phi_k(x) = e^{-\frac{(x-\mu_k)^2}{s^2}}$$

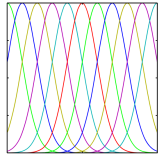
using 50 bases



● ● ●

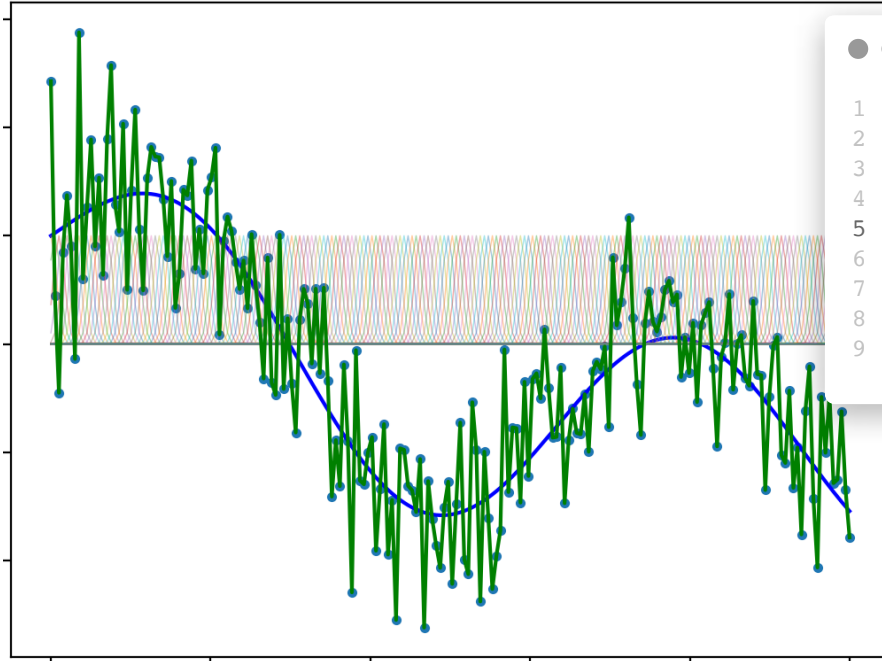
```
1 #x: N
2 #y: N
3 plt.plot(x, y, 'b.')
4 phi = lambda x,mu: np.exp(-(x-mu)**2)
5 mu = np.linspace(0,10,50) #50 Gaussian bases
6 Phi = phi(x[:,None], mu[None,:]) #N x 10
7 w = np.linalg.lstsq(Phi, y)[0]
8 yh = np.dot(Phi,w)
9 plt.plot(x, yh, 'g-')
```

# Example: Gaussian bases



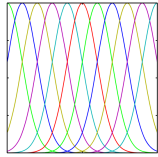
$$\phi_k(x) = e^{-\frac{(x-\mu_k)^2}{s^2}}$$

using 200, thinner bases ( $s=.1$ )



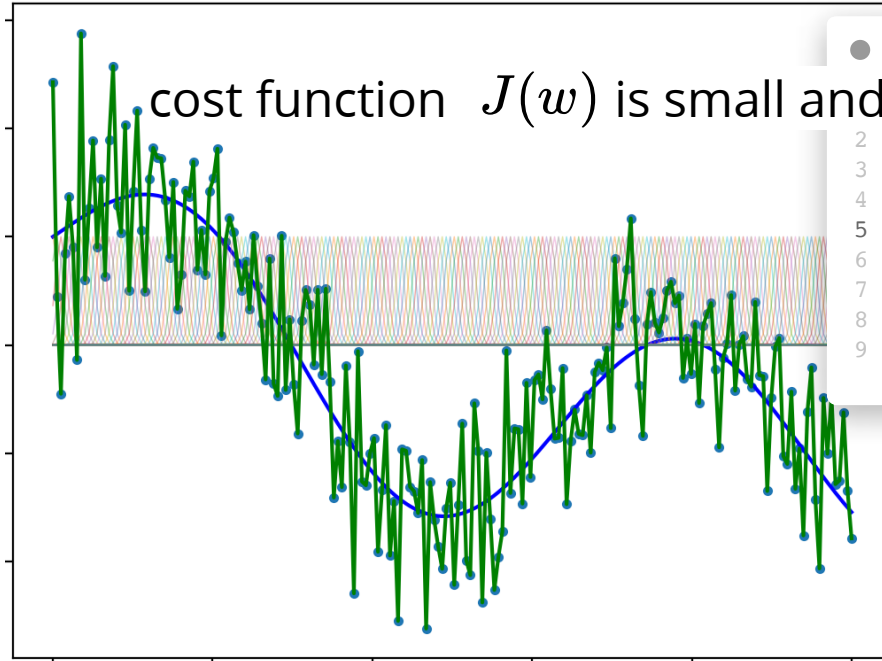
```
1 #x: N
2 #y: N
3 plt.plot(x, y, 'b.')
4 phi = lambda x,mu: np.exp(-((x-mu)/.1)**2)
5 mu = np.linspace(0,10,200) #200 Gaussians bases
6 Phi = phi(x[:,None], mu[None,:]) #N x 10
7 w = np.linalg.lstsq(Phi, y)[0]
8 yh = np.dot(Phi,w)
9 plt.plot(x, yh, 'g-')
```

# Example: Gaussian bases



$$\phi_k(x) = e^{-\frac{(x-\mu_k)^2}{s^2}}$$

using 200, thinner bases ( $s=.1$ )

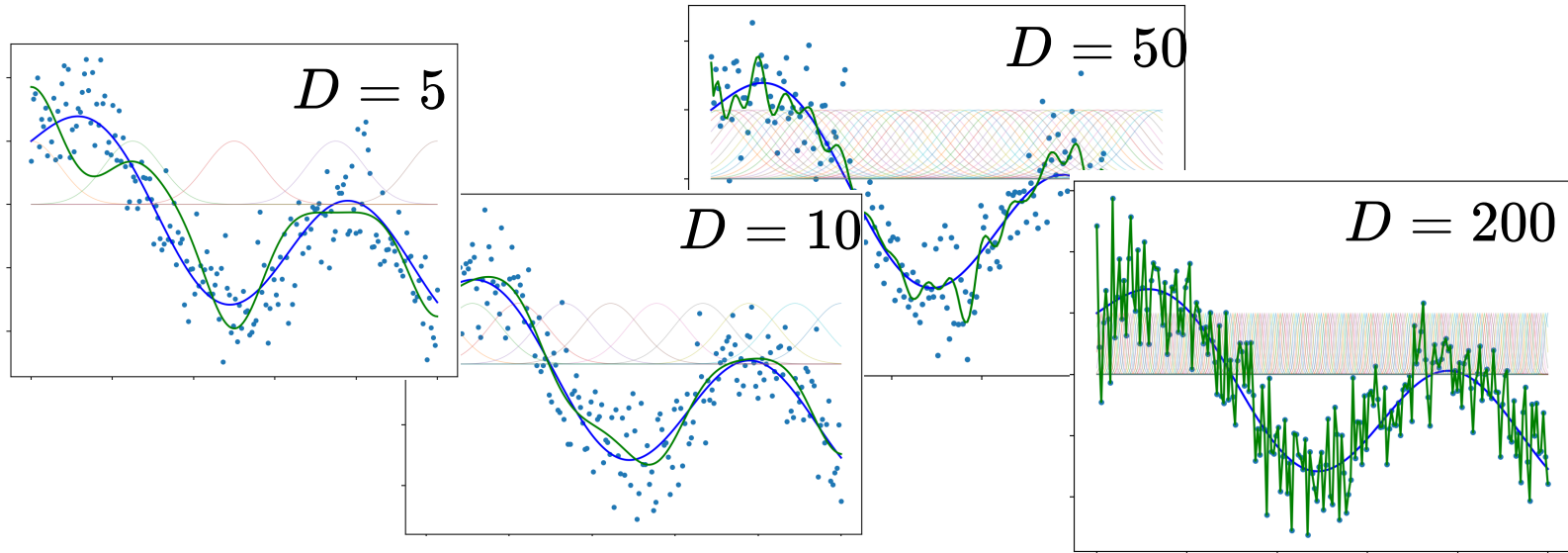


cost function  $J(w)$  is small and we have a "perfect" fit!

```
2 #y: N
3 plt.plot(x, y, 'b.')
4 phi = lambda x,mu: np.exp(-((x-mu)/.1)**2)
5 mu = np.linspace(0,10,200) #200 Gaussians bases
6 Phi = phi(x[:,None], mu[None,:]) #N x 10
7 w = np.linalg.lstsq(Phi, y)[0]
8 yh = np.dot(Phi,w)
9 plt.plot(x, yh, 'g-')
```



# Generalization

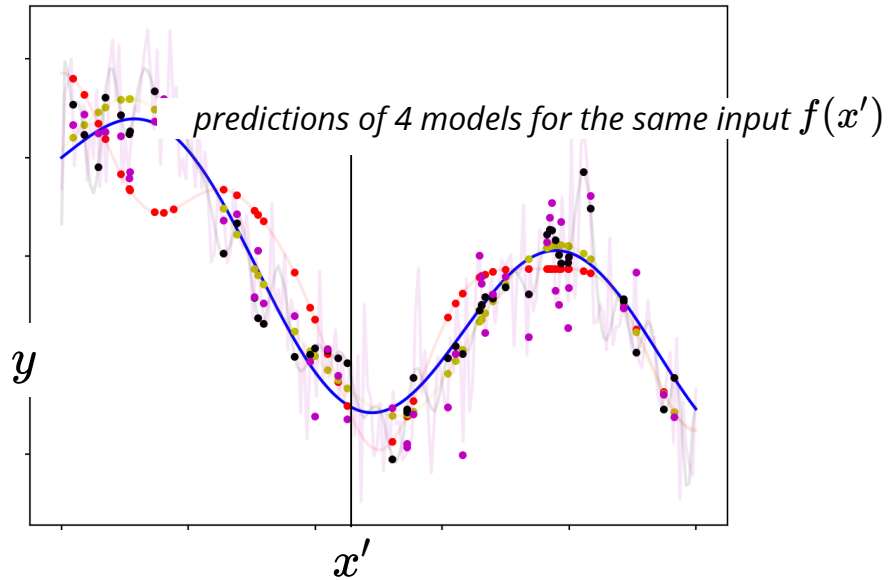


lower training error 

which one of these models performs better at **test time**?

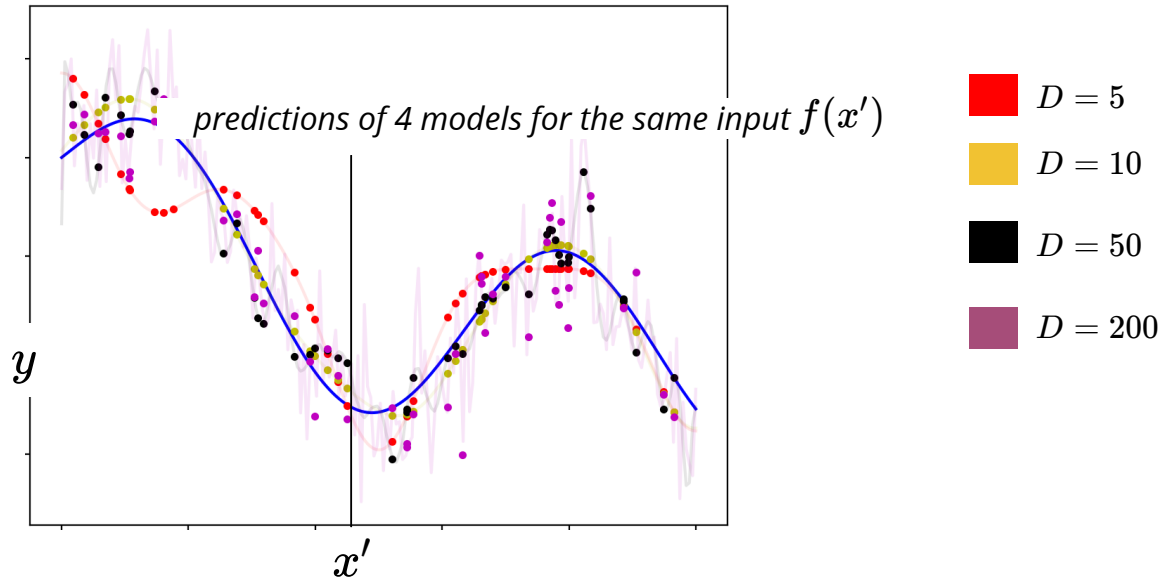
# Overfitting

which one of these models performs better at **test time**?



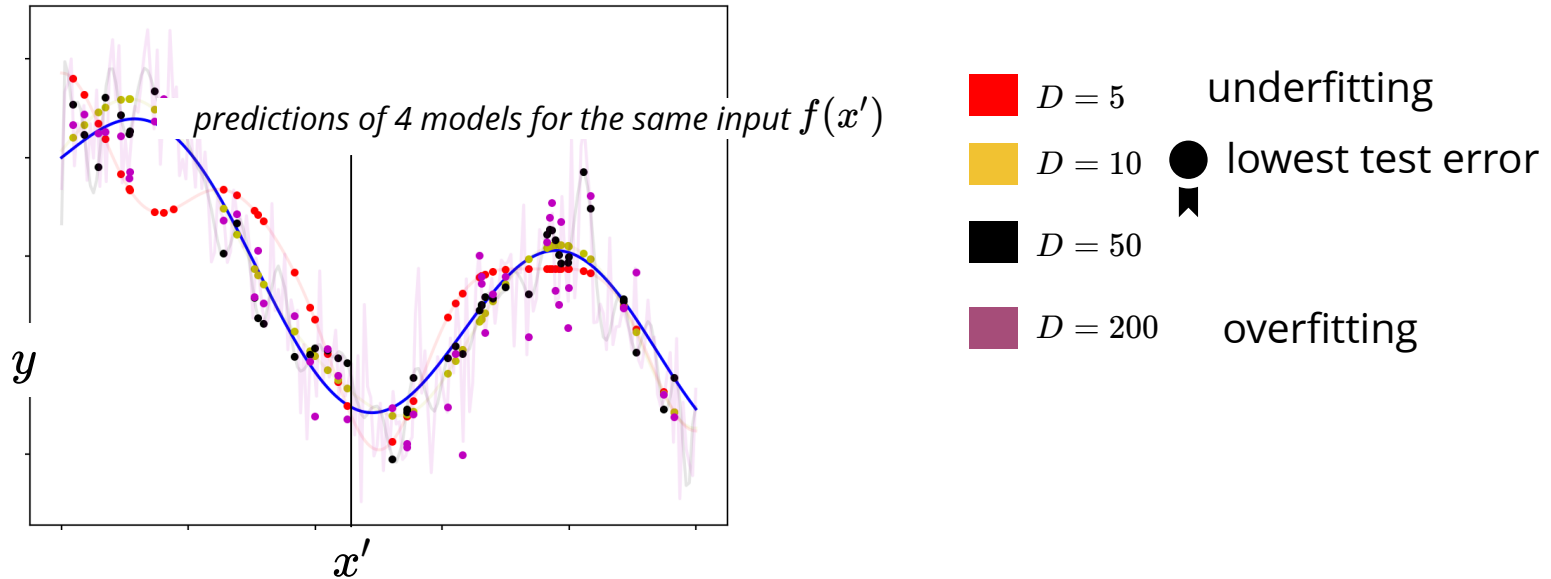
# Overfitting

which one of these models performs better at **test time**?



# Overfitting

which one of these models performs better at **test time**?



# Model selection

how to pick the model with lowest expected loss / test error?

**regularization** bound the test error by bounding

- training error
- model complexity

# Model selection

how to pick the model with lowest expected loss / test error?

**regularization** bound the test error by bounding

- training error
- model complexity

USE a **validation set** (and a separate test set for final assessment)



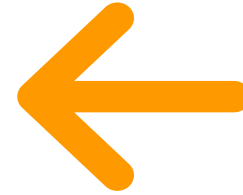
use for model selection    use for final model assessment

# Model selection

how to pick the model with lowest expected loss / test error?

**regularization** bound the test error by bounding

- training error
- model complexity



USE a **validation set** (and a separate test set for final assessment)



use for model selection

use for final model assessment

# An observation

when overfitting, we often see large weights



dashed lines are  $w_d \phi_d(x) \quad \forall d$

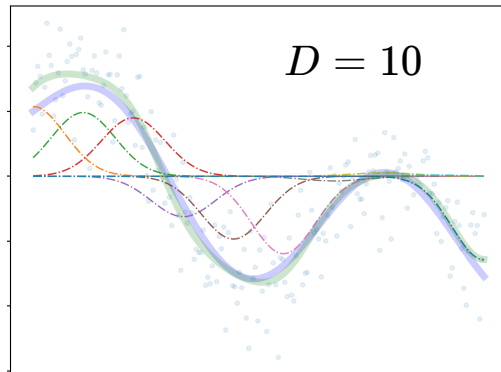


# An observation

when overfitting, we often see large weights



dashed lines are  $w_d \phi_d(x) \quad \forall d$

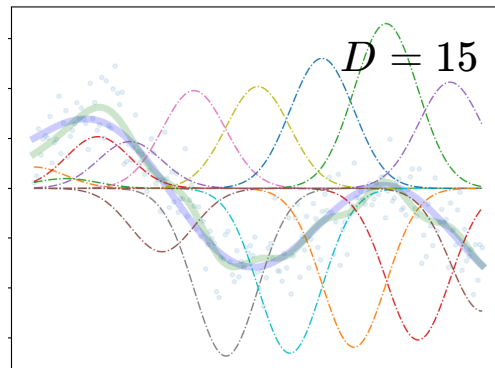
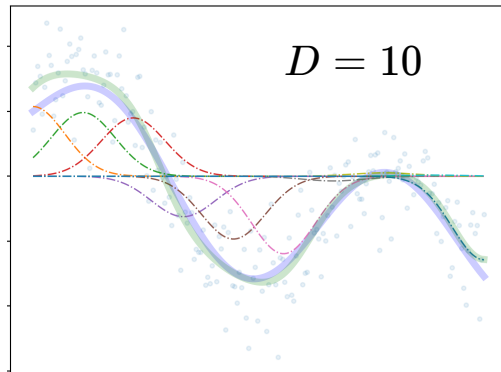


# An observation

when overfitting, we often see large weights



dashed lines are  $w_d \phi_d(x) \quad \forall d$

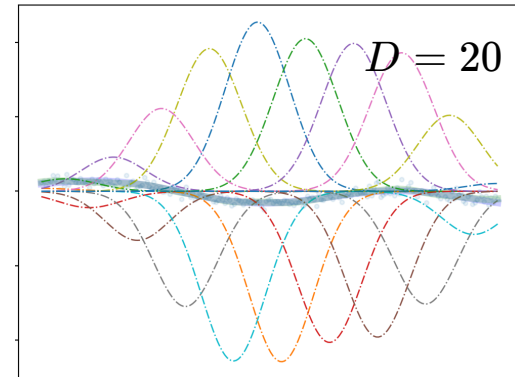
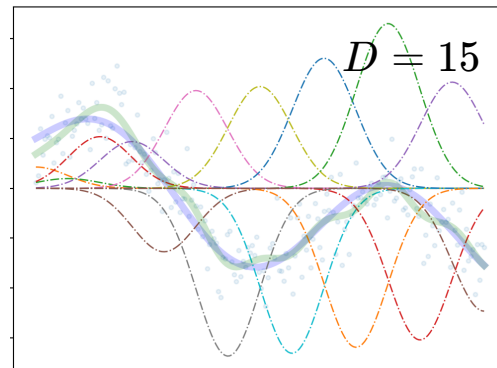
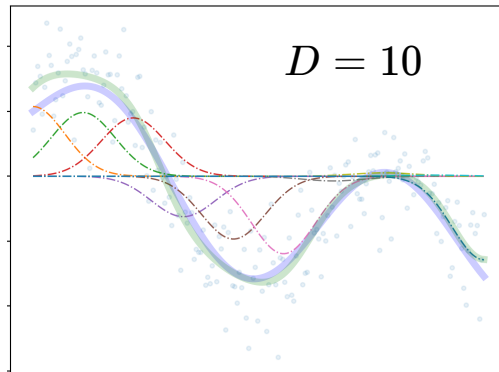


# An observation

when overfitting, we often see large weights



dashed lines are  $w_d \phi_d(x) \quad \forall d$

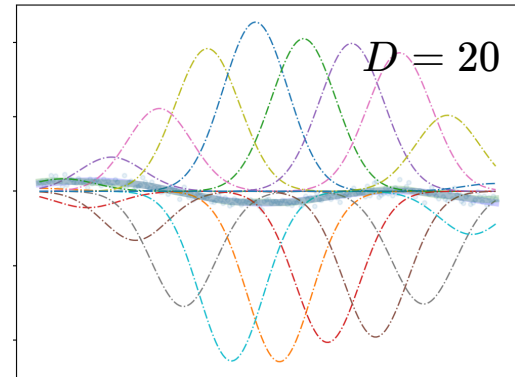
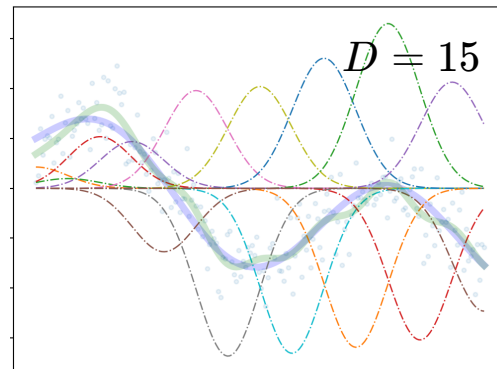
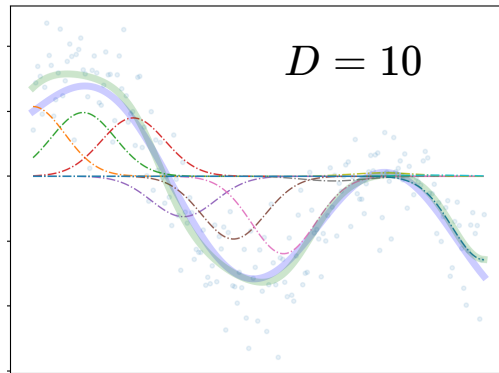


# An observation

when overfitting, we often see large weights



dashed lines are  $w_d \phi_d(x) \quad \forall d$



**idea:** penalize large parameter values

# Ridge regression

L2 regularized linear least squares regression:

$$J(w) = \frac{1}{2} \|Xw - y\|_2^2 + \frac{\lambda}{2} \|w\|_2^2$$

# Ridge regression

L2 regularized linear least squares regression:

$$J(w) = \frac{1}{2} \|Xw - y\|_2^2 + \frac{\lambda}{2} \|w\|_2^2$$

|  
sum of squared error  
 $\frac{1}{2} \sum_n (y^{(n)} - w^\top x)^2$

# Ridge regression

L2 regularized linear least squares regression:

$$J(w) = \frac{1}{2} \|Xw - y\|_2^2 + \frac{\lambda}{2} \|w\|_2^2$$

sum of squared error

$$\frac{1}{2} \sum_n (y^{(n)} - w^\top x)^2$$

(squared) L2 norm of w

$$w^T w = \sum_d w^2$$

# Ridge regression

L2 regularized linear least squares regression:

$$J(w) = \frac{1}{2} \|Xw - y\|_2^2 + \frac{\lambda}{2} \|w\|_2^2$$

$\downarrow$  |  
sum of squared error (squared) L2 norm of w

$$\frac{1}{2} \sum_n (y^{(n)} - w^\top x)^2 \qquad w^\top w = \sum_d w^2$$

regularization parameter  $\lambda > 0$  controls the strength of regularization



# Ridge regression

L2 regularized linear least squares regression:

$$J(w) = \frac{1}{2} \|Xw - y\|_2^2 + \frac{\lambda}{2} \|w\|_2^2$$

|

sum of squared error

$$\frac{1}{2} \sum_n (y^{(n)} - w^\top x)^2$$

|

(squared) L2 norm of w

$$w^\top w = \sum_d w^2$$

regularization parameter  $\lambda > 0$  controls the strength of regularization

a good practice is to not penalize the intercept  $\lambda(\|w\|_2^2 - w_0^2)$

# Ridge regression

we can set the derivative to zero  $J(w) = \frac{1}{2}(Xw - y)^\top(Xw - y) + \frac{\lambda}{2}w^\top w$   
 $\nabla J(w) = X^\top(Xw - y) + \lambda w = 0$

# Ridge regression

we can set the derivative to zero  $J(w) = \frac{1}{2}(Xw - y)^\top(Xw - y) + \frac{\lambda}{2}w^\top w$

$$\nabla J(w) = X^\top(Xw - y) + \lambda w = 0$$

when using gradient descent, this term reduces the weights at each step (**weight decay**)

# Ridge regression

we can set the derivative to zero  $J(w) = \frac{1}{2}(Xw - y)^\top(Xw - y) + \frac{\lambda}{2}w^\top w$

$$\nabla J(w) = X^\top(Xw - y) + \lambda w = 0$$

$$(X^\top X + \lambda I)w = X^\top y$$

when using gradient descent, this term reduces the weights at each step (**weight decay**)

# Ridge regression

we can set the derivative to zero  $J(w) = \frac{1}{2}(Xw - y)^\top(Xw - y) + \frac{\lambda}{2}w^\top w$

$$\nabla J(w) = X^\top(Xw - y) + \lambda w = 0$$

$$(X^\top X + \lambda I)w = X^\top y$$

when using gradient descent, this term reduces the weights at each step (**weight decay**)

$$w = (X^\top X + \lambda I)^{-1} X^\top y$$

# Ridge regression

we can set the derivative to zero  $J(w) = \frac{1}{2}(Xw - y)^\top(Xw - y) + \frac{\lambda}{2}w^\top w$

$$\nabla J(w) = X^\top(Xw - y) + \lambda w = 0$$

$$(X^\top X + \lambda I)w = X^\top y$$

when using gradient descent, this term reduces the weights at each step (**weight decay**)

$$w = (X^\top X + \lambda I)^{-1} X^\top y$$

the only part different due to regularization

# Ridge regression

we can set the derivative to zero  $J(w) = \frac{1}{2}(Xw - y)^\top(Xw - y) + \frac{\lambda}{2}w^\top w$

$$\nabla J(w) = X^\top(Xw - y) + \lambda w = 0$$

$$(X^\top X + \lambda I)w = X^\top y$$

when using gradient descent, this term reduces the weights at each step (**weight decay**)

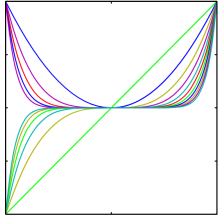
$$w = (X^\top X + \lambda I)^{-1} X^\top y$$

the only part different due to regularization

$\lambda I$  makes it invertible!

we can have linearly dependent features (e.g.,  $D > N$ )

the solution will be unique!



# Example: polynomial bases

polynomial bases

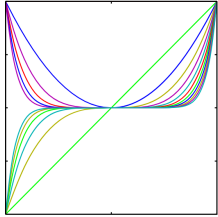
$$\phi_k(x) = x^k$$

Without regularization:

- using  $D=10$  we can perfectly fit the data (high test error)



# Example: polynomial bases

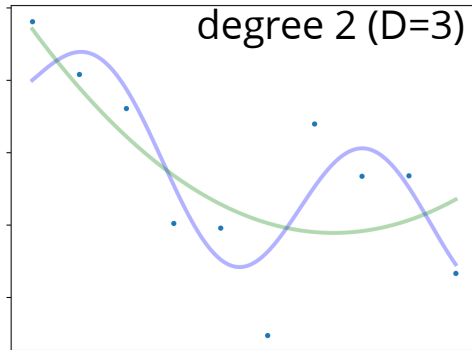


polynomial bases

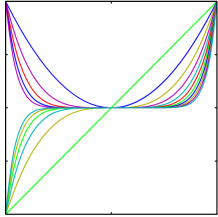
$$\phi_k(x) = x^k$$

Without regularization:

- using  $D=10$  we can perfectly fit the data (high test error)



# Example: polynomial bases

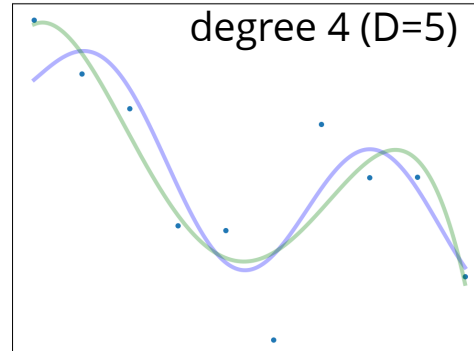
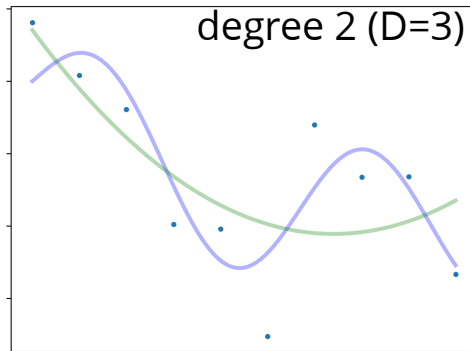


polynomial bases

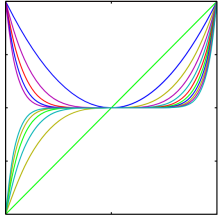
$$\phi_k(x) = x^k$$

Without regularization:

- using  $D=10$  we can perfectly fit the data (high test error)



# Example: polynomial bases

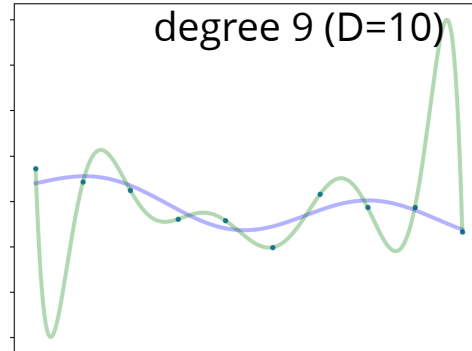
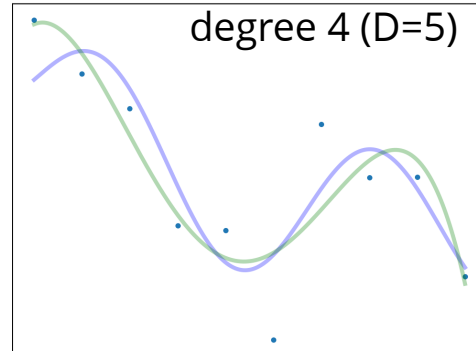
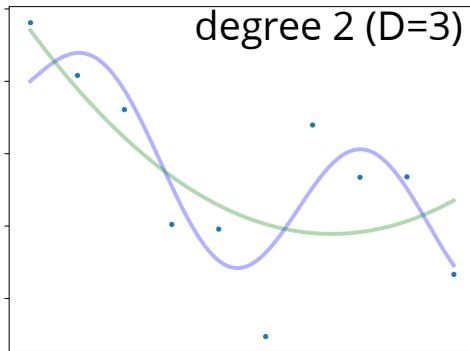


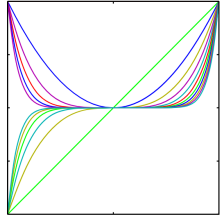
polynomial bases

$$\phi_k(x) = x^k$$

Without regularization:

- using  $D=10$  we can perfectly fit the data (high test error)





# Example: polynomial bases

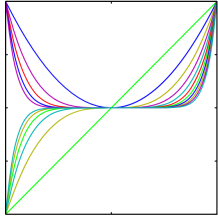
polynomial bases

$$\phi_k(x) = x^k$$

with regularization:

- fixed  $D=10$ , changing the amount of regularization

# Example: polynomial bases

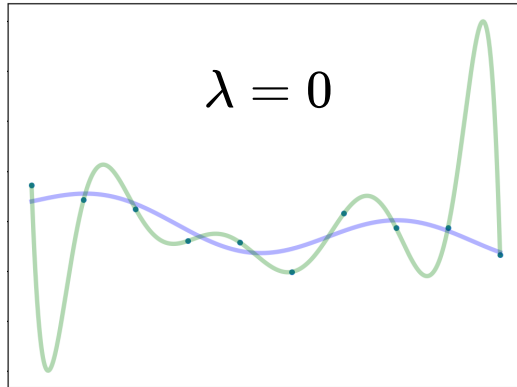


polynomial bases

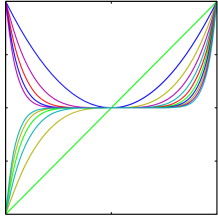
$$\phi_k(x) = x^k$$

with regularization:

- fixed  $D=10$ , changing the amount of regularization



# Example: polynomial bases

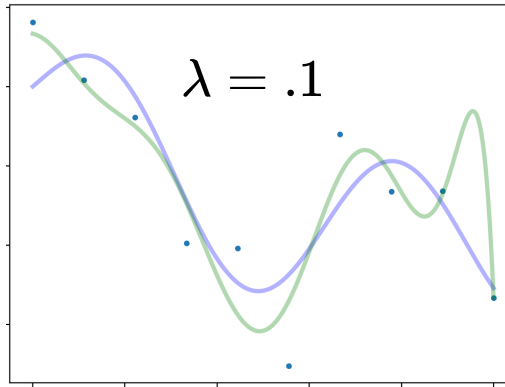
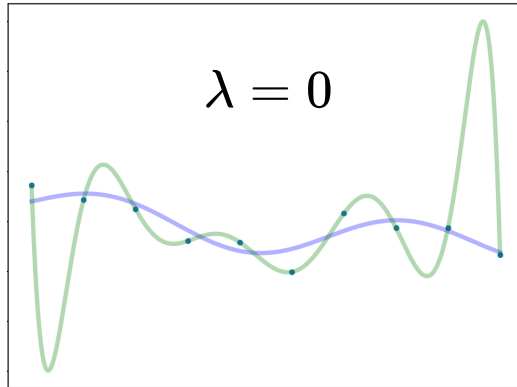


polynomial bases

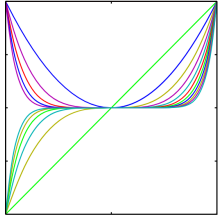
$$\phi_k(x) = x^k$$

with regularization:

- fixed  $D=10$ , changing the amount of regularization



# Example: polynomial bases

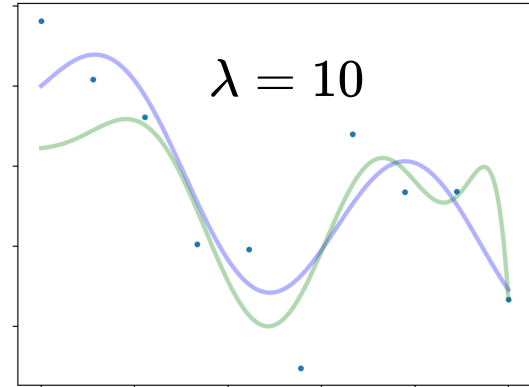
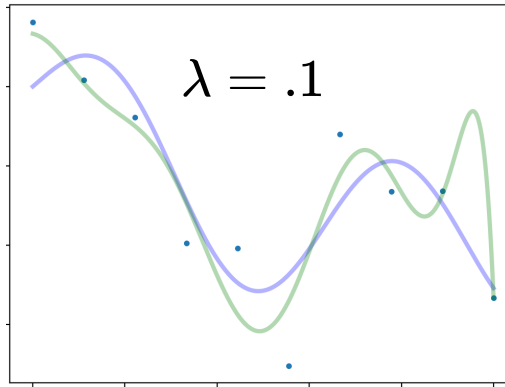
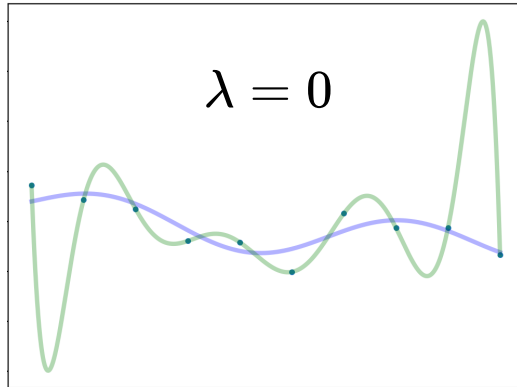


polynomial bases

$$\phi_k(x) = x^k$$

with regularization:

- fixed  $D=10$ , changing the amount of regularization



# Data normalization

what if we scale the input features, using different factors  $\tilde{x}^{(n)} = \gamma_d x^{(n)} \forall d, n$



# Data normalization

what if we scale the input features, using different factors  $\tilde{x}^{(n)} = \gamma_d x^{(n)} \forall d, n$

if we have **no regularization**:  $\tilde{w}_d = \frac{1}{\gamma_d} w_d \forall d$

everything remains the same because:  $\|Xw - y\|_2^2 = \|\tilde{X}\tilde{w} - y\|_2^2$

# Data normalization

what if we scale the input features, using different factors  $\tilde{x}^{(n)} = \gamma_d x^{(n)} \forall d, n$

if we have **no regularization**:  $\tilde{w}_d = \frac{1}{\gamma_d} w_d \forall d$

everything remains the same because:  $\|Xw - y\|_2^2 = \|\tilde{X}\tilde{w} - y\|_2^2$

**with regularization**:  $\|\tilde{w}\|_2 \neq \|w\|_2^2$  so the optimal **w** will be different!

# Data normalization

what if we scale the input features, using different factors  $\tilde{x}^{(n)} = \gamma_d x^{(n)} \forall d, n$

if we have **no regularization**:  $\tilde{w}_d = \frac{1}{\gamma_d} w_d \forall d$

everything remains the same because:  $\|Xw - y\|_2^2 = \|\tilde{X}\tilde{w} - y\|_2^2$

**with regularization**:  $\|\tilde{w}\|_2 \neq \|w\|_2^2$  so the optimal  $w$  will be different!

features of different mean and variance will be penalized differently

**normalization**  $\left\{ \begin{array}{l} \mu_d = \frac{1}{N} x_d^{(n)} \\ \sigma_d^2 = \frac{1}{N-1} (x_d^{(n)} - \mu_d)^2 \end{array} \right.$

makes sure all features have the same mean and variance  $x_d^{(n)} \leftarrow \frac{x_d^{(n)} - \mu_d}{\sigma_d}$

# Maximum likelihood

**previously:** linear regression & logistic regression maximize log-likelihood

# Maximum likelihood

**previously:** linear regression & logistic regression maximize log-likelihood

linear regression

$$\begin{aligned} w^* &= \arg \max_w p(y|w) \\ &= \arg \max_w \prod_{n=1}^N \mathcal{N}(y; \Phi w, \sigma^2) \\ &\equiv \arg \min \sum_n L_2(y^{(n)}, w^\top \phi(x^{(n)})) \end{aligned}$$

# Maximum likelihood

**previously:** linear regression & logistic regression maximize log-likelihood

linear regression

$$\begin{aligned}w^* &= \arg \max_w p(y|w) \\ &= \arg \max_w \prod_{n=1}^N \mathcal{N}(y; \Phi w, \sigma^2) \\ &\equiv \arg \min \sum_n L_2(y^{(n)}, w^\top \phi(x^{(n)}))\end{aligned}$$

logistic regression

$$\begin{aligned}w^* &= \arg \max_w p(y|x, w) \\ &= \arg \max_w \prod_{n=1}^N \text{Bernoulli}(y; \sigma(\Phi w)) \\ &\equiv \arg \min \sum_n L_{CE}(y^{(n)}, \sigma(w^\top \phi(x^n)))\end{aligned}$$

# Maximum likelihood

**previously:** linear regression & logistic regression maximize log-likelihood

linear regression

$$\begin{aligned}w^* &= \arg \max p(y|w) \\ &= \arg \max_w \prod_{n=1}^N \mathcal{N}(y; \Phi w, \sigma^2) \\ &\equiv \arg \min \sum_n L_2(y^{(n)}, w^\top \phi(x^{(n)}))\end{aligned}$$

logistic regression

$$\begin{aligned}w^* &= \arg \max p(y|x, w) \\ &= \arg \max_w \prod_{n=1}^N \text{Bernoulli}(y; \sigma(\Phi w)) \\ &\equiv \arg \min \sum_n L_{CE}(y^{(n)}, \sigma(w^\top \phi(x^n)))\end{aligned}$$

idea: maximize the posterior instead of likelihood

$$p(w|y) = \frac{p(w)p(y|w)}{p(y)}$$

# Maximum a Posteriori (MAP)

use the Bayes rule and find the parameters with max posterior prob.

$$p(w|y) = \frac{p(w)p(y|w)}{p(y)} \text{ the same for all choices of } w \text{ (ignore)}$$



# Maximum a Posteriori (MAP)

use the **Bayes rule** and find the parameters with max posterior prob.

MAP estimate  $p(w|y) = \frac{p(w)p(y|w)}{p(y)}$  the same for all choices of  $w$  (ignore)

$$w^* = \arg \max_w p(w)p(y|w)$$

# Maximum a Posteriori (MAP)

use the **Bayes rule** and find the parameters with max posterior prob.

MAP estimate 
$$p(w|y) = \frac{p(w)p(y|w)}{p(y)}$$
 the same for all choices of  $w$  (ignore)

$$w^* = \arg \max_w p(w)p(y|w)$$

$$\equiv \arg \max_w \log p(y|w) + \log p(w)$$

# Maximum a Posteriori (MAP)

use the **Bayes rule** and find the parameters with max posterior prob.

MAP estimate 
$$p(w|y) = \frac{p(w)p(y|w)}{p(y)}$$
 the same for all choices of  $w$  (ignore)

$$w^* = \arg \max_w p(w)p(y|w)$$

$$\equiv \arg \max_w \log p(y|w) + \log p(w)$$

likelihood: original objective

# Maximum a Posteriori (MAP)

use the **Bayes rule** and find the parameters with max posterior prob.

MAP estimate 
$$p(w|y) = \frac{p(w)p(y|w)}{p(y)}$$
 the same for all choices of  $w$  (ignore)

$$w^* = \arg \max_w p(w)p(y|w)$$

$$\equiv \arg \max_w \log p(y|w) + \log p(w)$$

likelihood: original objective    prior

# Maximum a Posteriori (MAP)

use the **Bayes rule** and find the parameters with max posterior prob.

MAP estimate 
$$p(w|y) = \frac{p(w)p(y|w)}{p(y)}$$
 the same for all choices of  $w$  (ignore)

$$w^* = \arg \max_w p(w)p(y|w)$$

$$\equiv \arg \max_w \log p(y|w) + \log p(w)$$

likelihood: original objective    prior

even better would be to estimate the posterior distribution  $p(w|y)$

- more on this later in the course!

# Gaussian prior

**Gaussian** likelihood and **Gaussian** prior

$$w^* = \arg \max_w p(w)p(y|w)$$

# Gaussian prior

**Gaussian** likelihood and **Gaussian** prior

$$w^* = \arg \max_w p(w)p(y|w) \equiv \arg \max_w \log p(y|w) + \log p(w)$$

# Gaussian prior

**Gaussian** likelihood and **Gaussian** prior

$$\begin{aligned}w^* &= \arg \max_w p(w)p(y|w) \equiv \arg \max_w \log p(y|w) + \log p(w) \\ &\equiv \arg \max_w \log \mathcal{N}(y|w^\top x, \sigma^2) + \sum_{d=1}^D \log \mathcal{N}(w_d, 0, \tau^2)\end{aligned}$$



# Gaussian prior

**Gaussian** likelihood and **Gaussian** prior

$$\begin{aligned} w^* &= \arg \max_w p(w)p(y|w) \equiv \arg \max_w \log p(y|w) + \log p(w) \\ &\equiv \arg \max_w \log \mathcal{N}(y|w^\top x, \sigma^2) + \sum_{d=1}^D \log \mathcal{N}(w_d, 0, \tau^2) \text{ assuming independent Gaussian} \\ &\hspace{15em} \text{(one per each weight)} \end{aligned}$$

# Gaussian prior

**Gaussian** likelihood and **Gaussian** prior

$$\begin{aligned}w^* &= \arg \max_w p(w)p(y|w) \equiv \arg \max_w \log p(y|w) + \log p(w) \\ &\equiv \arg \max_w \log \mathcal{N}(y|w^\top x, \sigma^2) + \sum_{d=1}^D \log \mathcal{N}(w_d, 0, \tau^2) \text{ assuming independent Gaussian} \\ &\quad \text{(one per each weight)} \\ &\equiv \arg \max_w \frac{-1}{2\sigma^2}(y - w^\top x)^2 - \sum_{d=1}^D \frac{1}{2\tau^2} w_d^2\end{aligned}$$

# Gaussian prior

**Gaussian** likelihood and **Gaussian** prior

$$\begin{aligned}w^* &= \arg \max_w p(w)p(y|w) \equiv \arg \max_w \log p(y|w) + \log p(w) \\ &\equiv \arg \max_w \log \mathcal{N}(y|w^\top x, \sigma^2) + \sum_{d=1}^D \log \mathcal{N}(w_d, 0, \tau^2) \text{ assuming independent Gaussian} \\ &\quad \text{(one per each weight)} \\ &\equiv \arg \max_w \frac{-1}{2\sigma^2}(y - w^\top x)^2 - \sum_{d=1}^D \frac{1}{2\tau^2} w_d^2 \equiv \arg \min_w \frac{1}{2}(y - w^\top x)^2 + \sum_{d=1}^D \frac{\sigma^2}{2\tau^2} w_d^2\end{aligned}$$

# Gaussian prior

**Gaussian** likelihood and **Gaussian** prior

$$\begin{aligned}w^* &= \arg \max_w p(w)p(y|w) \equiv \arg \max_w \log p(y|w) + \log p(w) \\ &\equiv \arg \max_w \log \mathcal{N}(y|w^\top x, \sigma^2) + \sum_{d=1}^D \log \mathcal{N}(w_d, 0, \tau^2) \text{ assuming independent Gaussian} \\ &\quad \text{(one per each weight)} \\ &\equiv \arg \max_w \frac{-1}{2\sigma^2}(y - w^\top x)^2 - \sum_{d=1}^D \frac{1}{2\tau^2} w_d^2 \equiv \arg \min_w \frac{1}{2}(y - w^\top x)^2 + \sum_{d=1}^D \frac{\sigma^2}{2\tau^2} w_d^2\end{aligned}$$

multiple data-points

$$\equiv \arg \min_w \frac{1}{2} \sum_n (y^{(n)} - w^\top x^{(n)})^2 + \sum_{d=1}^D \frac{\lambda}{2} w_d^2$$

# Gaussian prior

**Gaussian** likelihood and **Gaussian** prior

$$\begin{aligned}w^* &= \arg \max_w p(w)p(y|w) \equiv \arg \max_w \log p(y|w) + \log p(w) \\ &\equiv \arg \max_w \log \mathcal{N}(y|w^\top x, \sigma^2) + \sum_{d=1}^D \log \mathcal{N}(w_d, 0, \tau^2) \text{ assuming independent Gaussian} \\ &\quad \text{(one per each weight)} \\ &\equiv \arg \max_w \frac{-1}{2\sigma^2}(y - w^\top x)^2 - \sum_{d=1}^D \frac{1}{2\tau^2} w_d^2 \equiv \arg \min_w \frac{1}{2}(y - w^\top x)^2 + \sum_{d=1}^D \frac{\sigma^2}{2\tau^2} w_d^2\end{aligned}$$

multiple data-points

$$\begin{aligned}&\lambda = \frac{\sigma^2}{\tau^2} \\ &\quad \downarrow \\ &\equiv \arg \min_w \frac{1}{2} \sum_n (y^{(n)} - w^\top x^{(n)})^2 + \sum_{d=1}^D \frac{\lambda}{2} w_d^2 \quad \text{L2 regularization}\end{aligned}$$

# Gaussian prior

**Gaussian** likelihood and **Gaussian** prior

$$\begin{aligned}w^* &= \arg \max_w p(w)p(y|w) \equiv \arg \max_w \log p(y|w) + \log p(w) \\ &\equiv \arg \max_w \log \mathcal{N}(y|w^\top x, \sigma^2) + \sum_{d=1}^D \log \mathcal{N}(w_d, 0, \tau^2) \text{ assuming independent Gaussian} \\ &\quad \text{(one per each weight)} \\ &\equiv \arg \max_w \frac{-1}{2\sigma^2}(y - w^\top x)^2 - \sum_{d=1}^D \frac{1}{2\tau^2} w_d^2 \equiv \arg \min_w \frac{1}{2}(y - w^\top x)^2 + \sum_{d=1}^D \frac{\sigma^2}{2\tau^2} w_d^2\end{aligned}$$

multiple data-points

$$\begin{aligned}&\lambda = \frac{\sigma^2}{\tau^2} \\ &\quad \downarrow \\ &\equiv \arg \min_w \frac{1}{2} \sum_n (y^{(n)} - w^\top x^{(n)})^2 + \sum_{d=1}^D \frac{\lambda}{2} w_d^2 \quad \text{L2 regularization}\end{aligned}$$

**L2- regularization** is assuming a **Gaussian prior** on weights

the same is true for logistic regression (or any other cost function)

# Laplace prior

another notable choice of prior is the Laplace distribution

image:<https://stats.stackexchange.com/questions/177210/why-is-laplace-prior-producing-sparse-solutions>

# Laplace prior

another notable choice of prior is the Laplace distribution

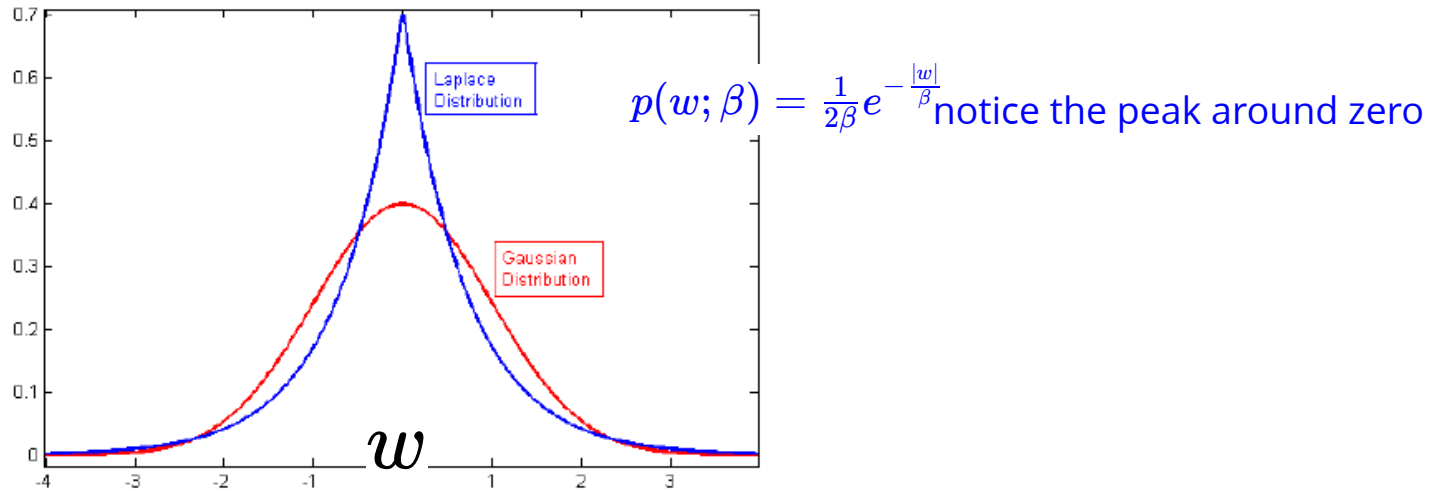


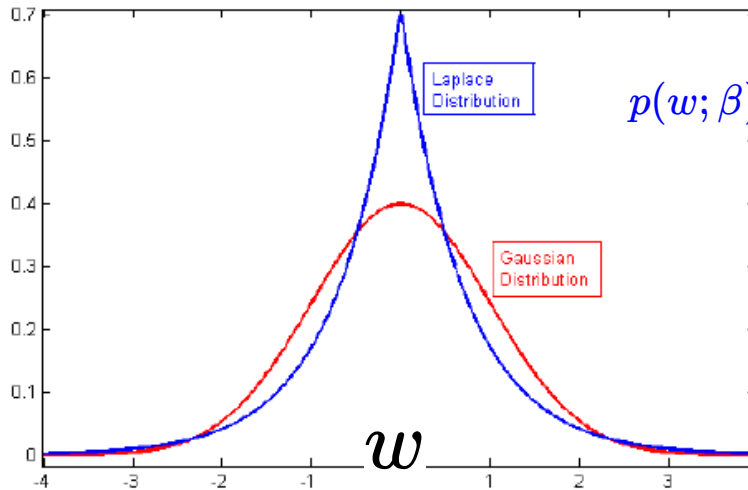
image:<https://stats.stackexchange.com/questions/177210/why-is-laplace-prior-producing-sparse-solutions>



# Laplace prior

another notable choice of prior is the Laplace distribution

minimizing negative log-likelihood  $\rightarrow -\sum_d \log p(w_d) = \sum_d \frac{1}{2\beta} |w_d|$



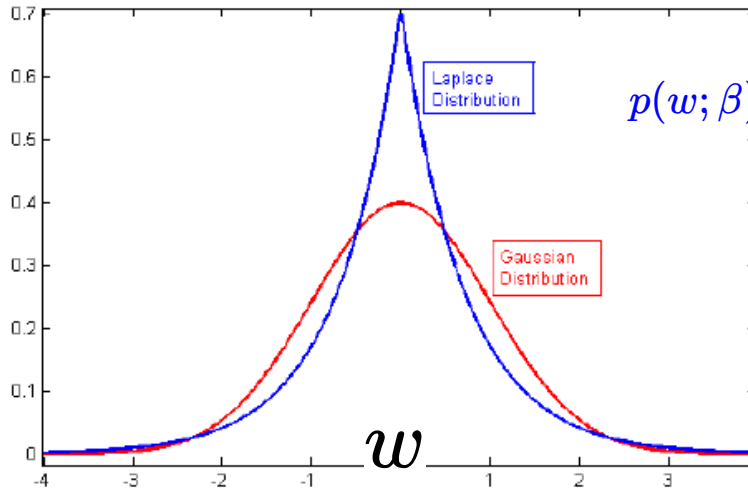
$$p(w; \beta) = \frac{1}{2\beta} e^{-\frac{|w|}{\beta}} \text{ notice the peak around zero}$$

image:<https://stats.stackexchange.com/questions/177210/why-is-laplace-prior-producing-sparse-solutions>

# Laplace prior

another notable choice of prior is the Laplace distribution

minimizing negative log-likelihood  $\rightarrow -\sum_d \log p(w_d) = \sum_d \frac{1}{2\beta} |w_d| = \frac{1}{2\beta} \|w\|_1$   
L1 norm of  $w$



$$p(w; \beta) = \frac{1}{2\beta} e^{-\frac{|w|}{\beta}} \text{ notice the peak around zero}$$

image: <https://stats.stackexchange.com/questions/177210/why-is-laplace-prior-producing-sparse-solutions>

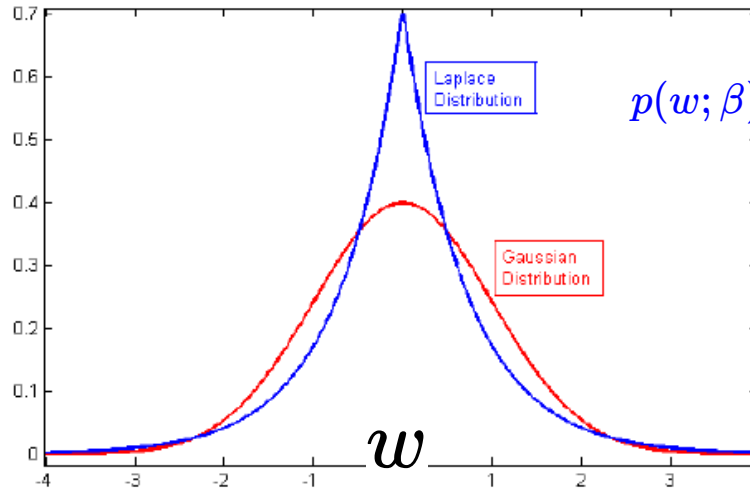
# Laplace prior

another notable choice of prior is the Laplace distribution

minimizing negative log-likelihood  $\rightarrow -\sum_d \log p(w_d) = \sum_d \frac{1}{2\beta} |w_d| = \frac{1}{2\beta} \|w\|_1$   
L1 norm of  $w$

L1 regularization:  $J(w) \leftarrow J(w) + \lambda \|w\|_1$  also called **lasso**

(least absolute shrinkage and selection operator)

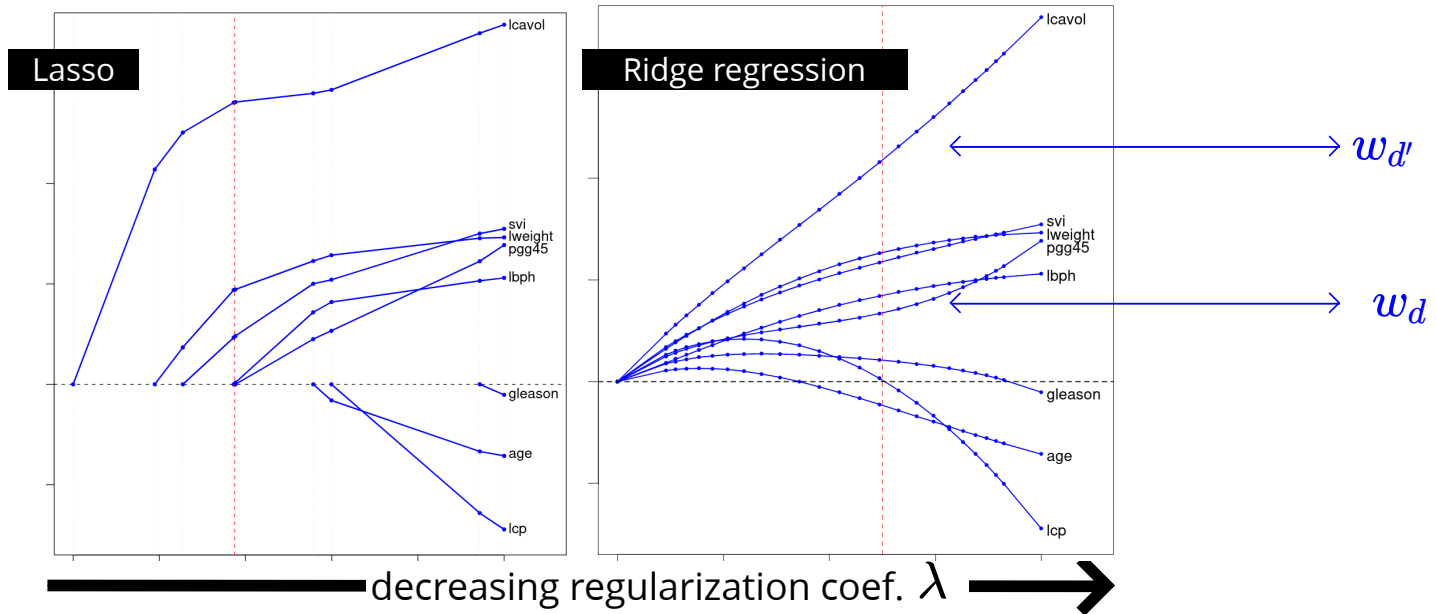


$$p(w; \beta) = \frac{1}{2\beta} e^{-\frac{|w|}{\beta}} \text{ notice the peak around zero}$$

image: <https://stats.stackexchange.com/questions/177210/why-is-laplace-prior-producing-sparse-solutions>

# $L_1$ vs $L_2$ regularization

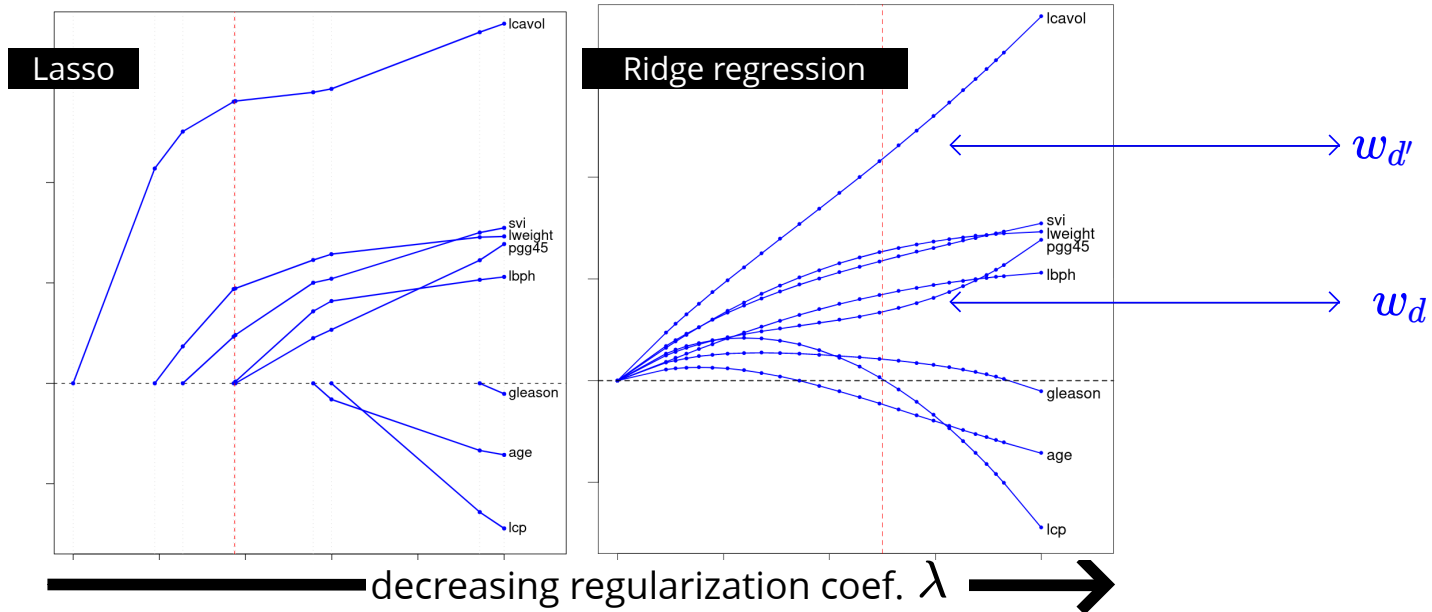
regularization path shows how  $\{w_d\}$  change as we change  $\lambda$



# $L_1$ vs $L_2$ regularization

regularization path shows how  $\{w_d\}$  change as we change  $\lambda$

Lasso produces sparse weights (many are zero, rather than small)

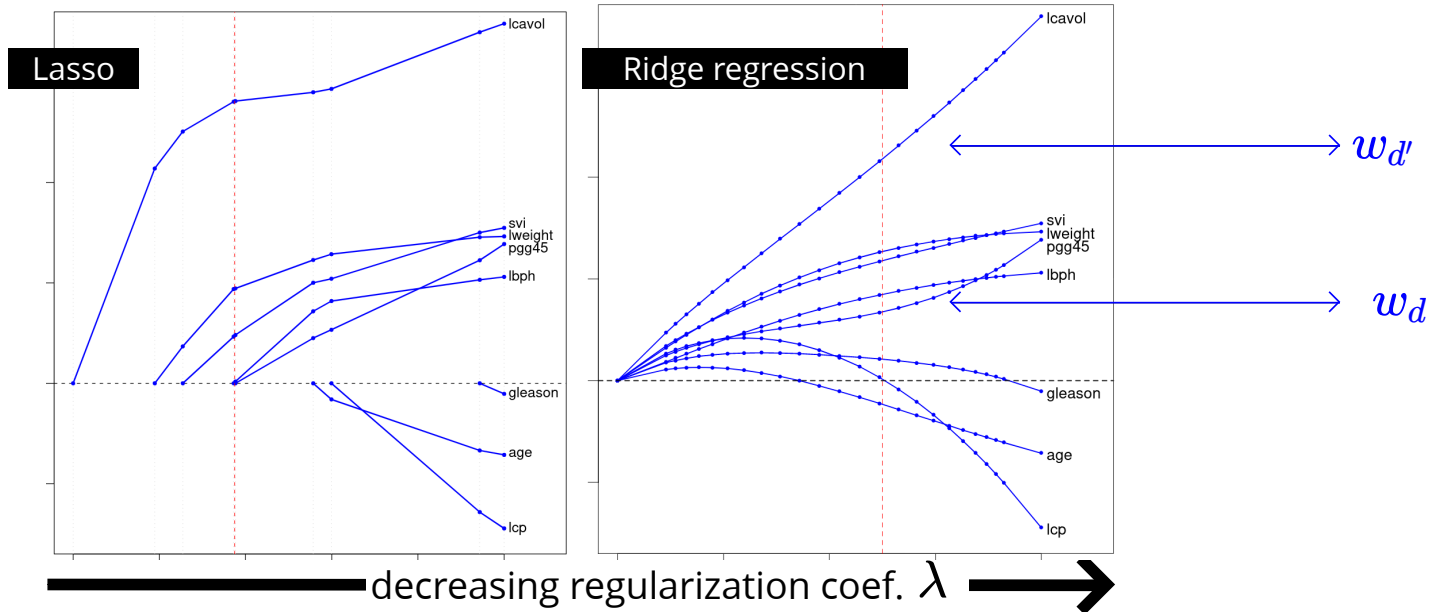


# $L_1$ vs $L_2$ regularization

regularization path shows how  $\{w_d\}$  change as we change  $\lambda$

Lasso produces sparse weights (many are zero, rather than small)

red-line is the optimal  $\lambda$  from cross-validation



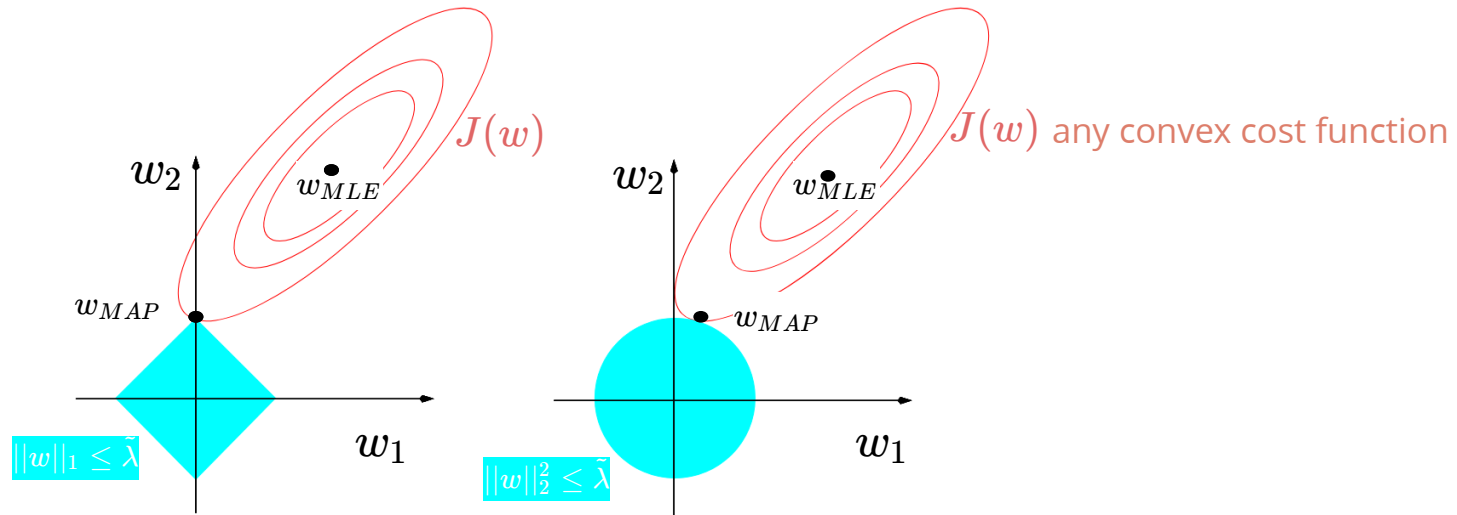
# $L_1$ vs $L_2$ regularization

$\min_w J(w) + \lambda \|w\|_p^p$   
is equivalent to  $\min_w J(w)$  subject to  $\|w\|_p^p \leq \tilde{\lambda}$  for an appropriate choice of  $\tilde{\lambda}$

# $L_1$ vs $L_2$ regularization

$$\min_w J(w) + \lambda \|w\|_p^p$$

is equivalent to  $\min_w J(w)$  subject to  $\|w\|_p^p \leq \tilde{\lambda}$  for an appropriate choice of  $\tilde{\lambda}$   
figures below show the constraint and the isocontours of  $J(w)$





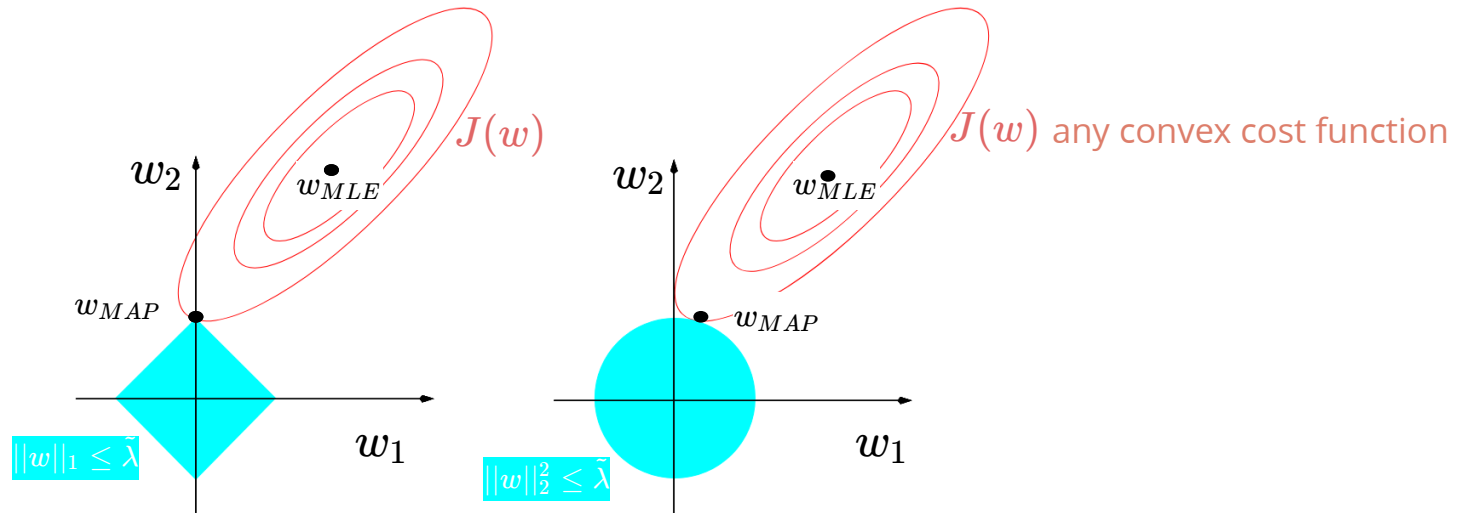
# $L_1$ vs $L_2$ regularization

$$\min_w J(w) + \lambda \|w\|_p^p$$

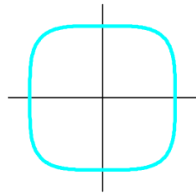
is equivalent to  $\min_w J(w)$  subject to  $\|w\|_p^p \leq \tilde{\lambda}$  for an appropriate choice of  $\tilde{\lambda}$

figures below show the constraint and the isocontours of  $J(w)$

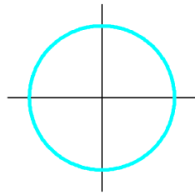
optimal solution with L1-regularization is more likely to have zero components



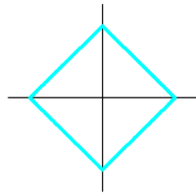
# Subset selection



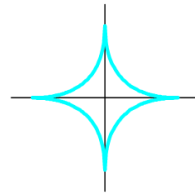
$$\sum_d w_d^4$$



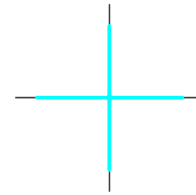
$$\sum_d w_d^2$$



$$\sum_d |w_d|$$



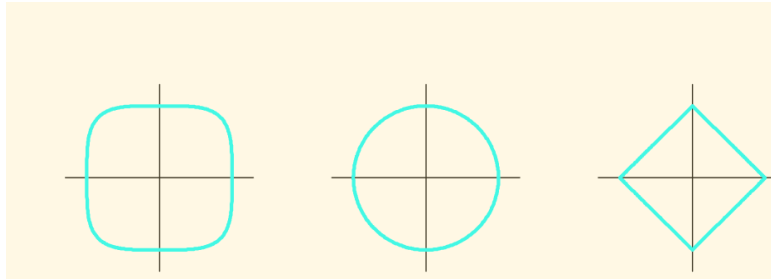
$$\sum_d |w_d|^{\frac{1}{2}}$$



$$\sum_d |w_d|^{\frac{1}{10}}$$

# Subset selection

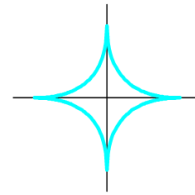
p-norms with  $p \geq 1$  are convex (easier to optimize)



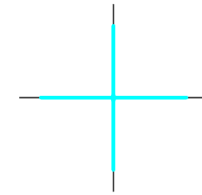
$$\sum_d w_d^4$$

$$\sum_d w_d^2$$

$$\sum_d |w_d|$$



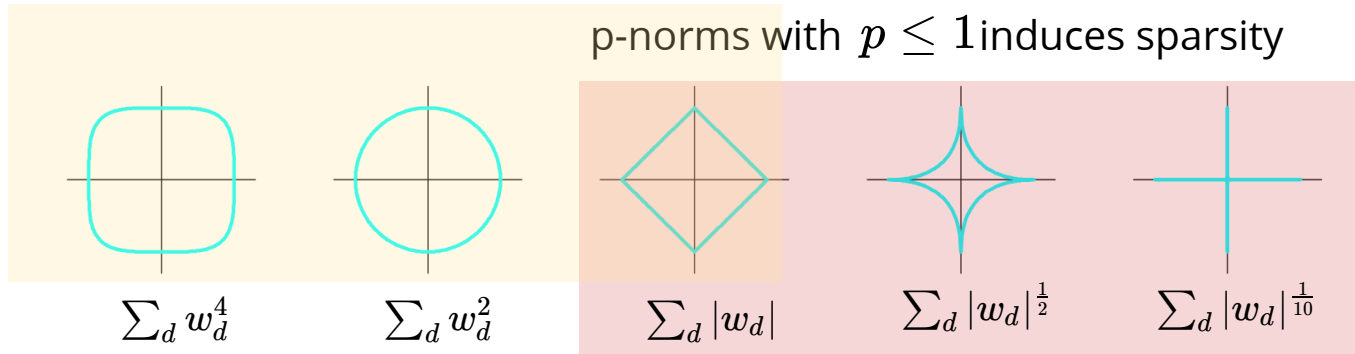
$$\sum_d |w_d|^{\frac{1}{2}}$$



$$\sum_d |w_d|^{\frac{1}{10}}$$

# Subset selection

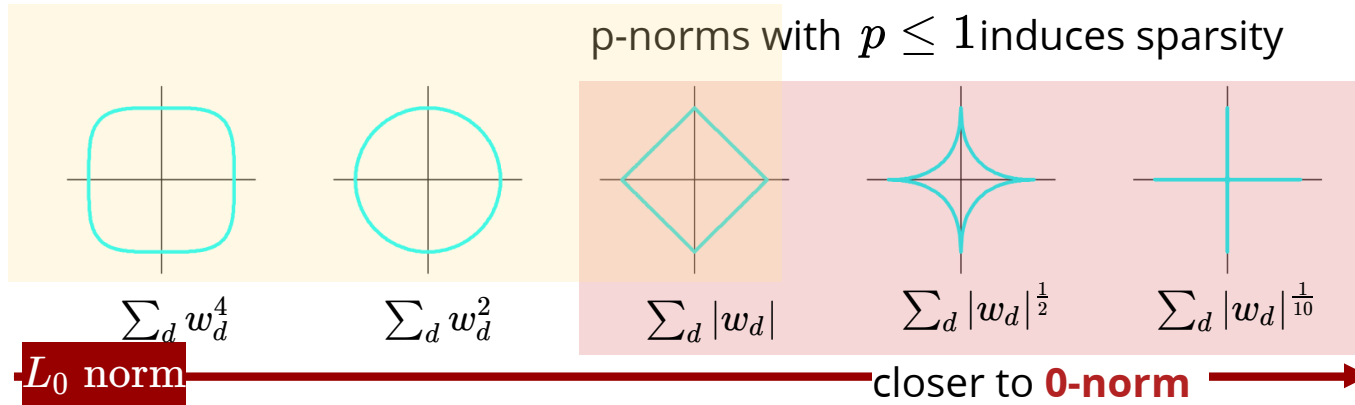
p-norms with  $p \geq 1$  are convex (easier to optimize)



# Subset selection

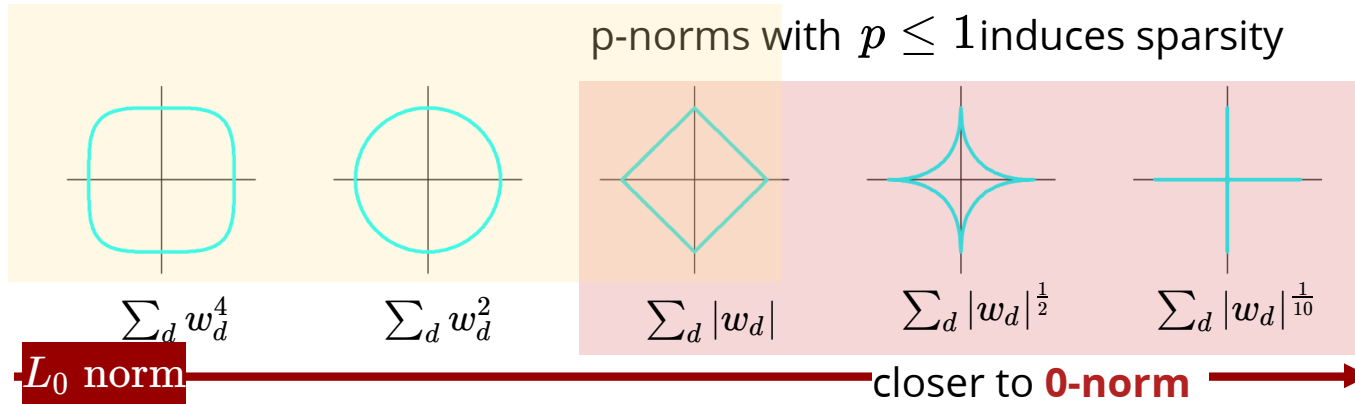
p-norms with  $p \geq 1$  are convex (easier to optimize)

p-norms with  $p \leq 1$  induces sparsity



# Subset selection

p-norms with  $p \geq 1$  are convex (easier to optimize)

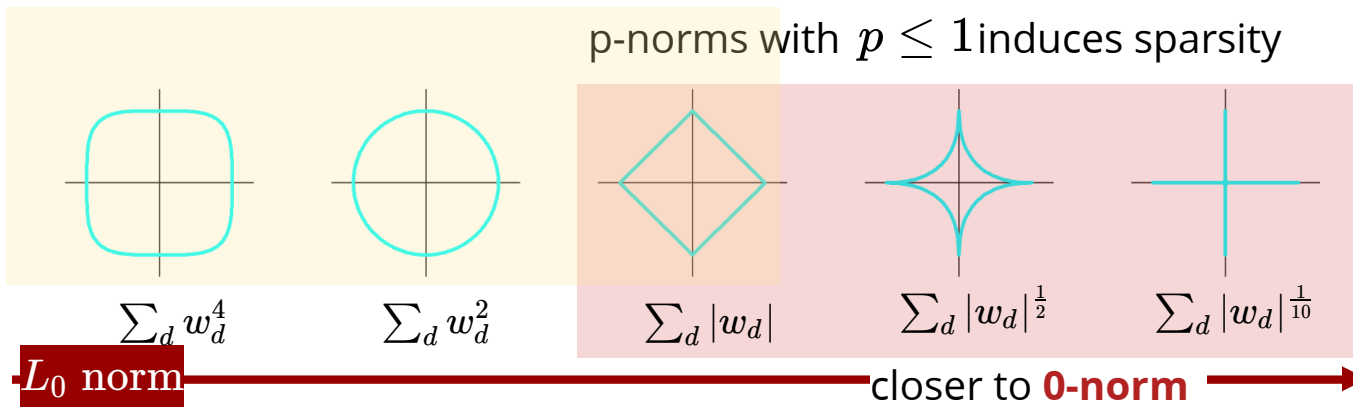


penalizes the **number of** non-zero features

$$J(w) + \lambda \|w\|_0 = J(w) + \lambda \sum_d \mathbb{I}(w_d \neq 0)$$

# Subset selection

p-norms with  $p \geq 1$  are convex (easier to optimize)



penalizes the **number of** non-zero features

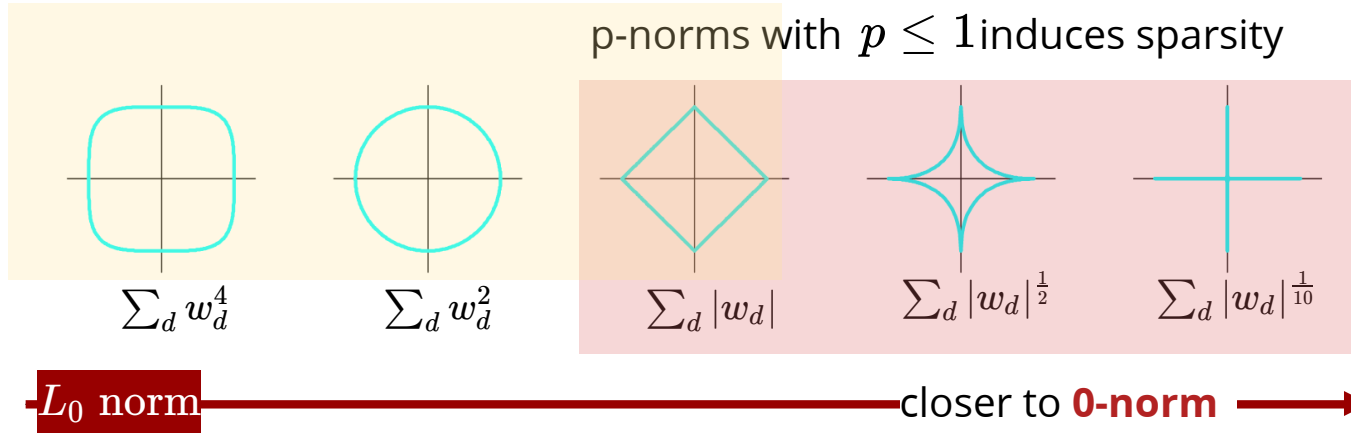
$$J(w) + \lambda \|w\|_0 = J(w) + \lambda \sum_d \mathbb{I}(w_d \neq 0)$$

a penalty of  $\lambda$  for each feature

performs feature selection

# Subset selection

p-norms with  $p \geq 1$  are convex (easier to optimize)



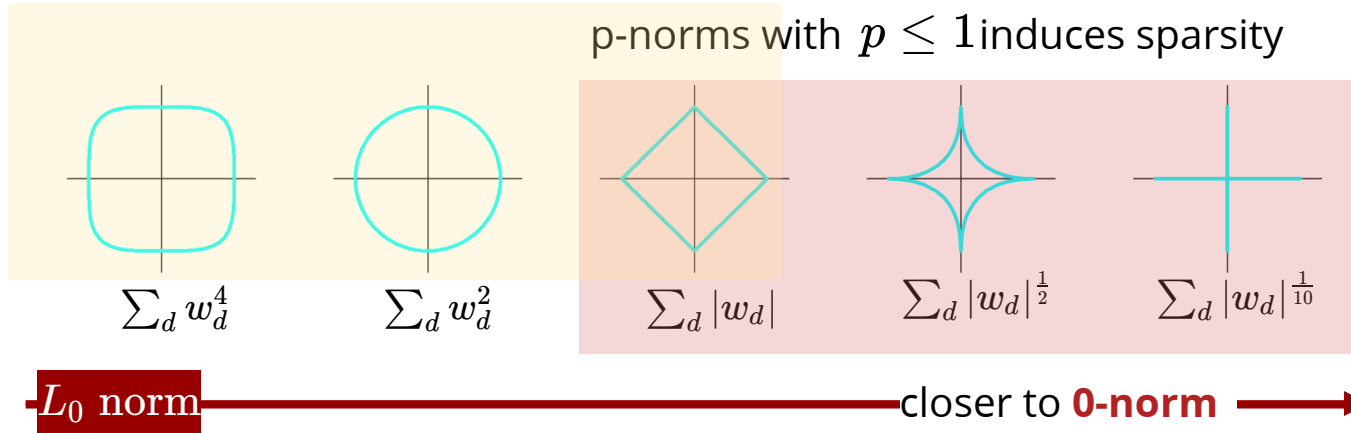
optimizing this is a difficult *combinatorial problem*:

- search over all  $2^D$  subsets



# Subset selection

p-norms with  $p \geq 1$  are convex (easier to optimize)



optimizing this is a difficult *combinatorial problem*:

- search over all  $2^D$  subsets

L1 regularization is a viable alternative to L0 regularization

# Bias-variance decomposition

for L2 loss

# Bias-variance decomposition

for L2 loss

assume a true distribution  $p(x, y)$

# Bias-variance decomposition

for L2 loss

assume a true distribution  $p(x, y)$

the regression function is  $f(x) = \mathbb{E}_p[y|x]$

# Bias-variance decomposition

for L2 loss

assume a true distribution  $p(x, y)$

the regression function is  $f(x) = \mathbb{E}_p[y|x]$

assume that a dataset  $\mathcal{D} = \{(x^{(n)}, y^{(n)})\}_n$  is sampled from  $p(x, y)$

# Bias-variance decomposition

for L2 loss

assume a true distribution  $p(x, y)$

the regression function is  $f(x) = \mathbb{E}_p[y|x]$

assume that a dataset  $\mathcal{D} = \{(x^{(n)}, y^{(n)})\}_n$  is sampled from  $p(x, y)$

let  $\hat{f}_{\mathcal{D}}$  be our model based on the dataset

# Bias-variance decomposition

for L2 loss

assume a true distribution  $p(x, y)$

the regression function is  $f(x) = \mathbb{E}_p[y|x]$

assume that a dataset  $\mathcal{D} = \{(x^{(n)}, y^{(n)})\}_n$  is sampled from  $p(x, y)$

let  $\hat{f}_{\mathcal{D}}$  be our model based on the dataset

what we care about is the **expected loss (aka risk)**

$$\mathbb{E}[(\hat{f}_{\mathcal{D}}(x) - y)^2]$$

all **blue** items are **random variables**

# Bias-variance decomposition

for L2 loss

what we care about is the *expected loss (aka risk)*

$$\mathbb{E}[(\hat{f}_{\mathcal{D}}(x) - y)^2]$$



# Bias-variance decomposition

for L2 loss

what we care about is the *expected loss (aka risk)*

$$\mathbb{E}[(\hat{f}_{\mathcal{D}}(x) - \underset{\substack{| \\ f(x) + \epsilon}}{y})^2]$$

# Bias-variance decomposition

for L2 loss

what we care about is the *expected loss (aka risk)*

$$\mathbb{E}[(\hat{f}_{\mathcal{D}}(x) - y)^2]$$

$\hat{f}_{\mathcal{D}}(x) + \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)]$  add and subtract a term

$f(x) + \epsilon$

# Bias-variance decomposition

for L2 loss

what we care about is the *expected loss (aka risk)*

$$\mathbb{E}[(\hat{f}_{\mathcal{D}}(x) - y)^2] = \mathbb{E}[(\hat{f}_{\mathcal{D}}(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - y + \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2]$$

$\hat{f}_{\mathcal{D}}(x)$        $f(x) + \epsilon$

$\hat{f}_{\mathcal{D}}(x) + \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)]$  add and subtract a term

# Bias-variance decomposition

for L2 loss

what we care about is the *expected loss (aka risk)*

$$\mathbb{E}[(\hat{f}_{\mathcal{D}}(x) - y)^2] = \mathbb{E}[(\hat{f}_{\mathcal{D}}(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - y + \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2]$$

$\hat{f}_{\mathcal{D}}(x)$        $f(x) + \epsilon$

|                      |

$\hat{f}_{\mathcal{D}}(x) + \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)]$       add and subtract a term

the remaining terms evaluate to zero (check for yourself!)

$$= \mathbb{E}[(\hat{f}_{\mathcal{D}}(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2] + \mathbb{E}[(f(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2] + \mathbb{E}[\epsilon^2]$$

# Bias-variance decomposition

for L2 loss

what we care about is the *expected loss (aka risk)*

$$\mathbb{E}[(\hat{f}_{\mathcal{D}}(x) - y)^2] = \mathbb{E}[(\hat{f}_{\mathcal{D}}(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - y + \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2]$$

$\hat{f}_{\mathcal{D}}(x)$        $f(x) + \epsilon$

|                      |

$\hat{f}_{\mathcal{D}}(x) + \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)]$       add and subtract a term

the remaining terms evaluate to zero (check for yourself!)

$$= \mathbb{E}[(\hat{f}_{\mathcal{D}}(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2] + \mathbb{E}[(f(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2] + \mathbb{E}[\epsilon^2]$$

variance    bias    unavoidable noise error

# Bias-variance decomposition

for L2 loss

the expected loss is decomposed to:

image: P. Domingos' posted article

# Bias-variance decomposition

for L2 loss

the expected loss is decomposed to:

$$\mathbb{E}[(f(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2]$$

**bias:** how average over all datasets differs from the regression function

image: P. Domingos' posted article

# Bias-variance decomposition

for L2 loss

the expected loss is decomposed to:

$$\mathbb{E}[(f(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2]$$

**bias:** how average over all datasets differs from the regression function

$$\mathbb{E}[(\hat{f}_{\mathcal{D}}(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2]$$

**variance:** how change of dataset affects the prediction

image: P. Domingos' posted article



# Bias-variance decomposition

for L2 loss

the expected loss is decomposed to:

$$\mathbb{E}[(f(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2]$$

**bias:** how average over all datasets differs from the regression function

$$\mathbb{E}[(\hat{f}_{\mathcal{D}}(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2]$$

**variance:** how change of dataset affects the prediction

$$\mathbb{E}[\epsilon^2]$$

**noise error:** the error even if we used the true model  $f(x)$

image: P. Domingos' posted article

# Bias-variance decomposition

for L2 loss

the expected loss is decomposed to:

$$\mathbb{E}[(f(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2]$$

**bias:** how average over all datasets differs from the regression function

$$\mathbb{E}[(\hat{f}_{\mathcal{D}}(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2]$$

**variance:** how change of dataset affects the prediction

$$\mathbb{E}[\epsilon^2]$$

**noise error:** the error even if we used the true model  $f(x)$

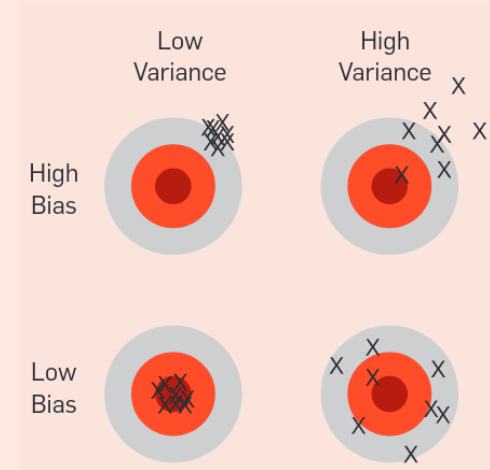


image: P. Domingos' posted article

# Bias-variance decomposition

for L2 loss

the expected loss is decomposed to:

$$\mathbb{E}[(f(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2]$$

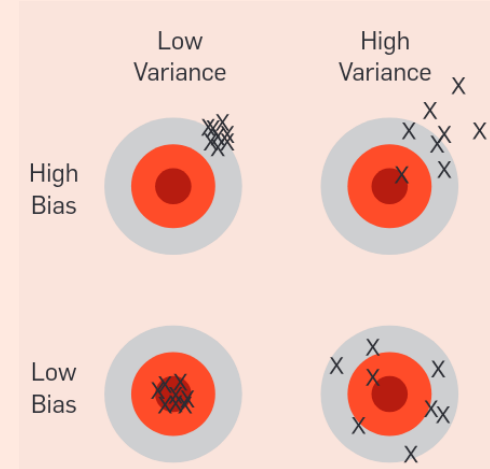
**bias:** how average over all datasets differs from the regression function

$$\mathbb{E}[(\hat{f}_{\mathcal{D}}(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2]$$

**variance:** how change of dataset affects the prediction

$$\mathbb{E}[\epsilon^2]$$

**noise error:** the error even if we used the true model  $f(x)$

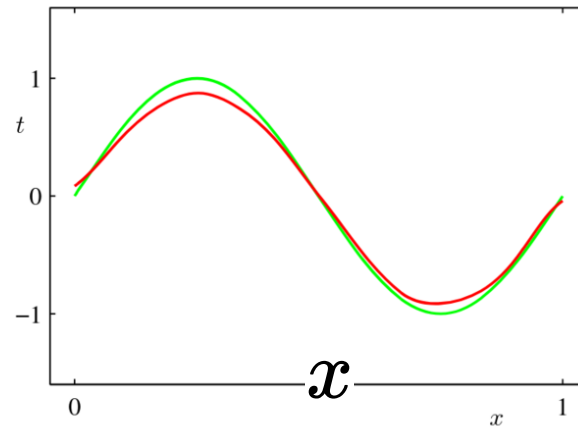
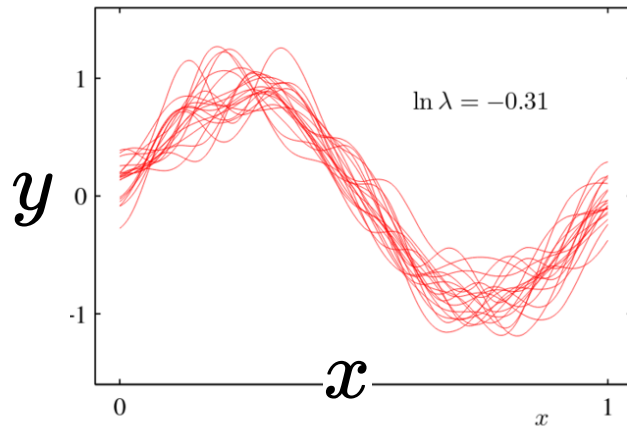


different models vary in their trade off between error due to bias and variance

- simple models: often more biased
- complex models: often have more variance

image: P. Domingos' posted article

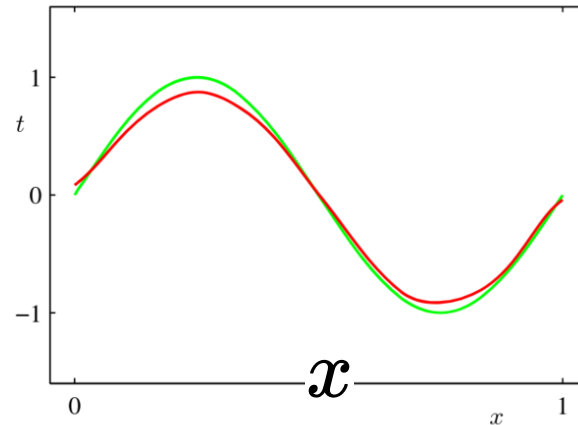
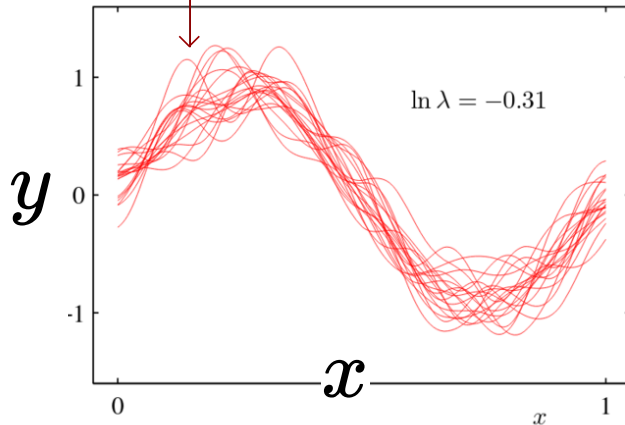
# Example: bias vs. variance



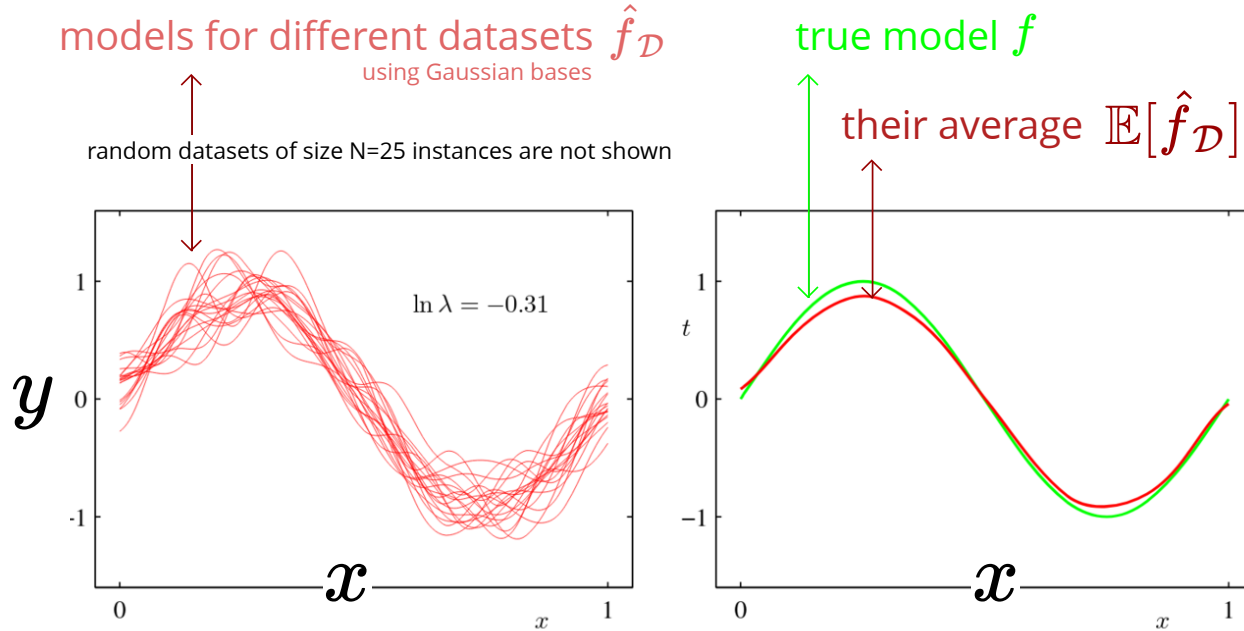
# Example: bias vs. variance

models for different datasets  $\hat{f}_{\mathcal{D}}$   
using Gaussian bases

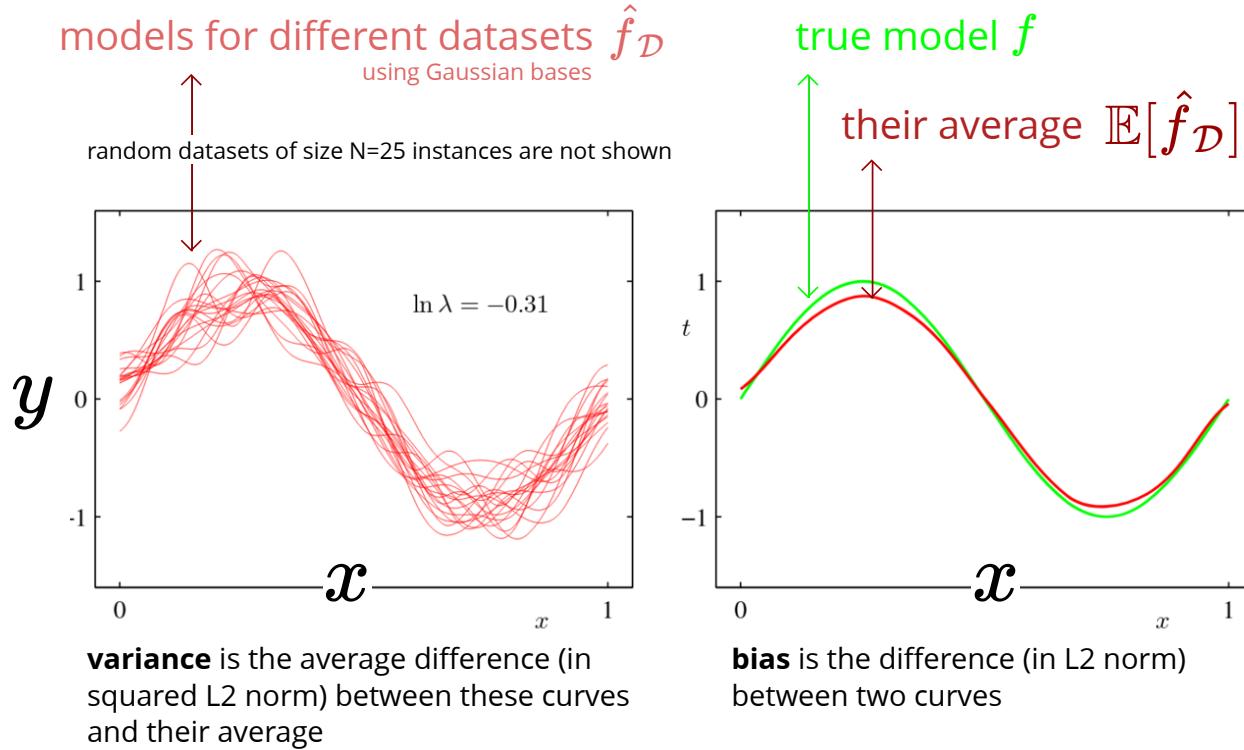
random datasets of size  $N=25$  instances are not shown



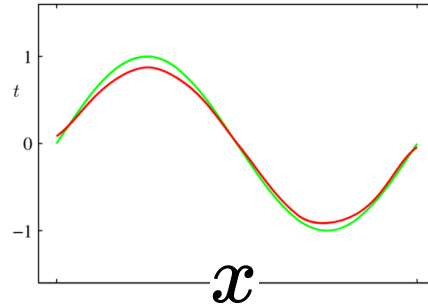
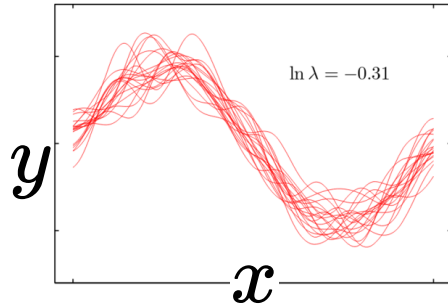
# Example: bias vs. variance



# Example: bias vs. variance

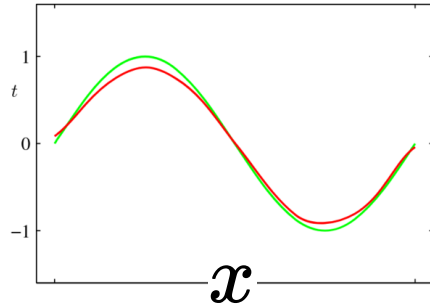
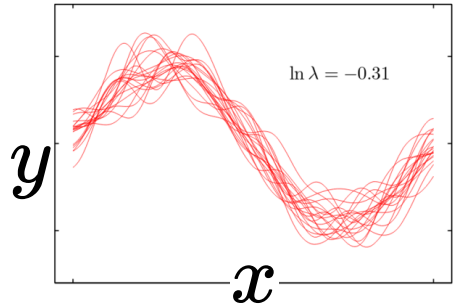


# Example: bias vs. variance

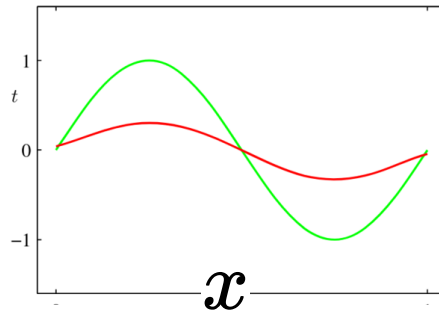
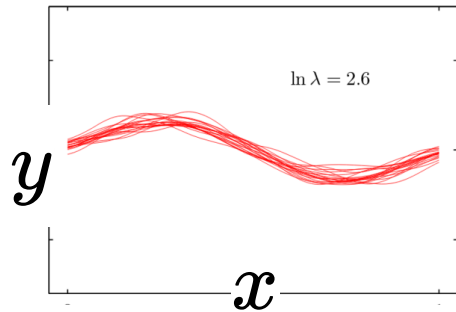




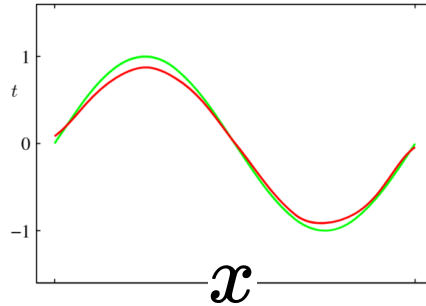
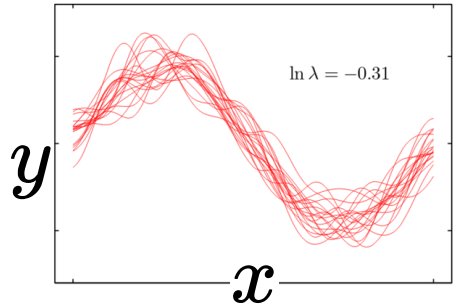
# Example: bias vs. variance



using larger regularization penalty: higher bias - lower variance



# Example: bias vs. variance

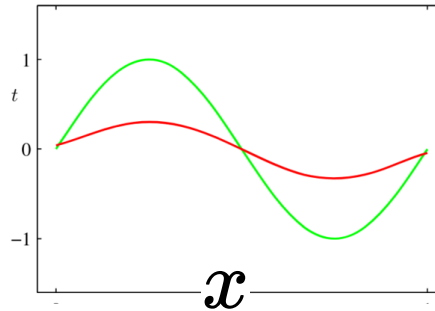
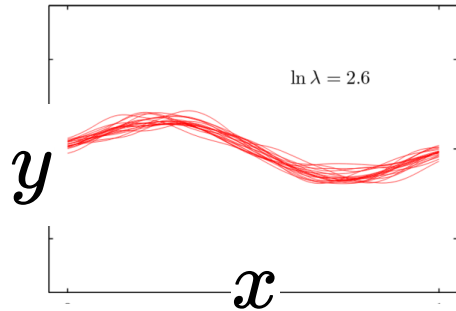


side note

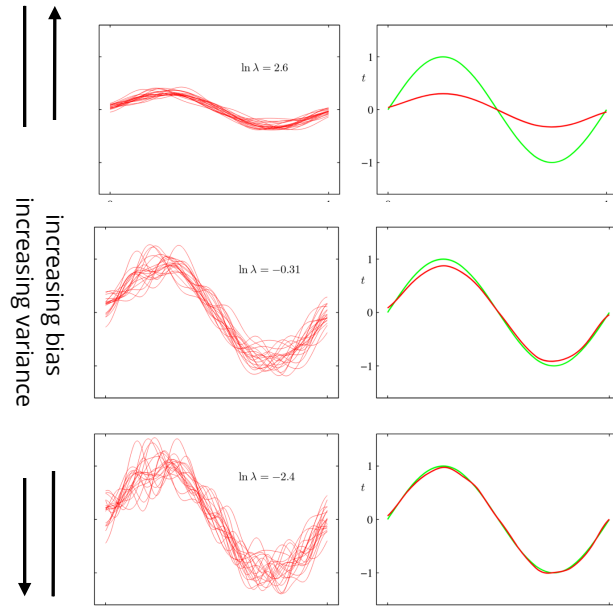
the average fit is very good, despite high variance

**model averaging:** uses "average" prediction of expressive models to prevent overfitting

using larger regularization penalty: higher bias - lower variance

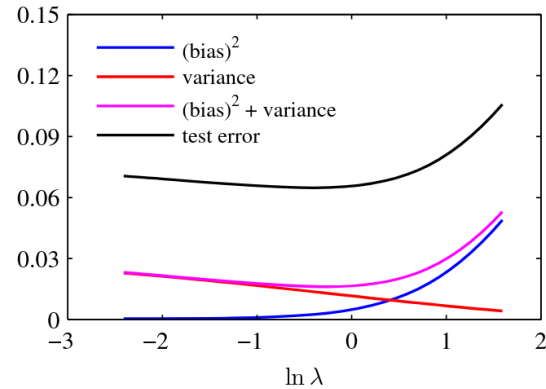
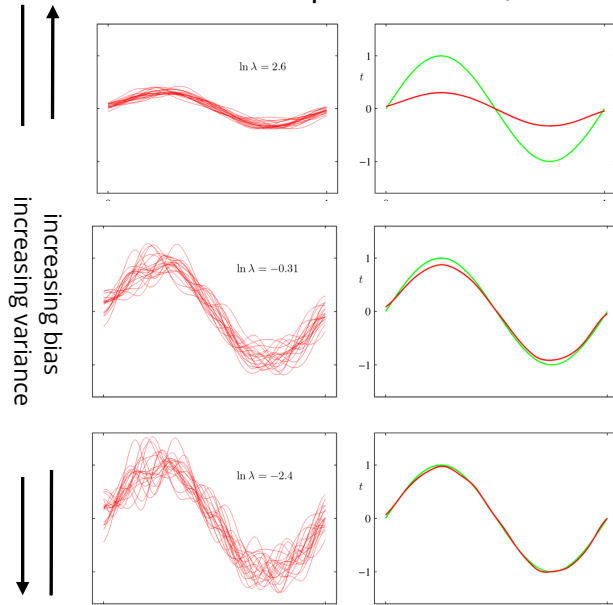


# Example: bias vs. variance



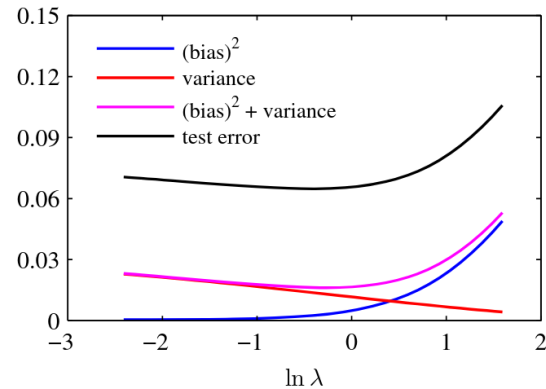
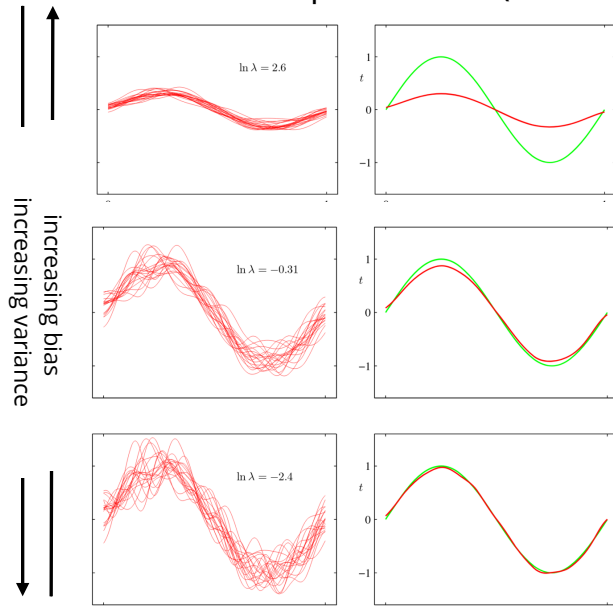
# Example: bias vs. variance

the lowest expected loss (test error) is somewhere between the two extremes



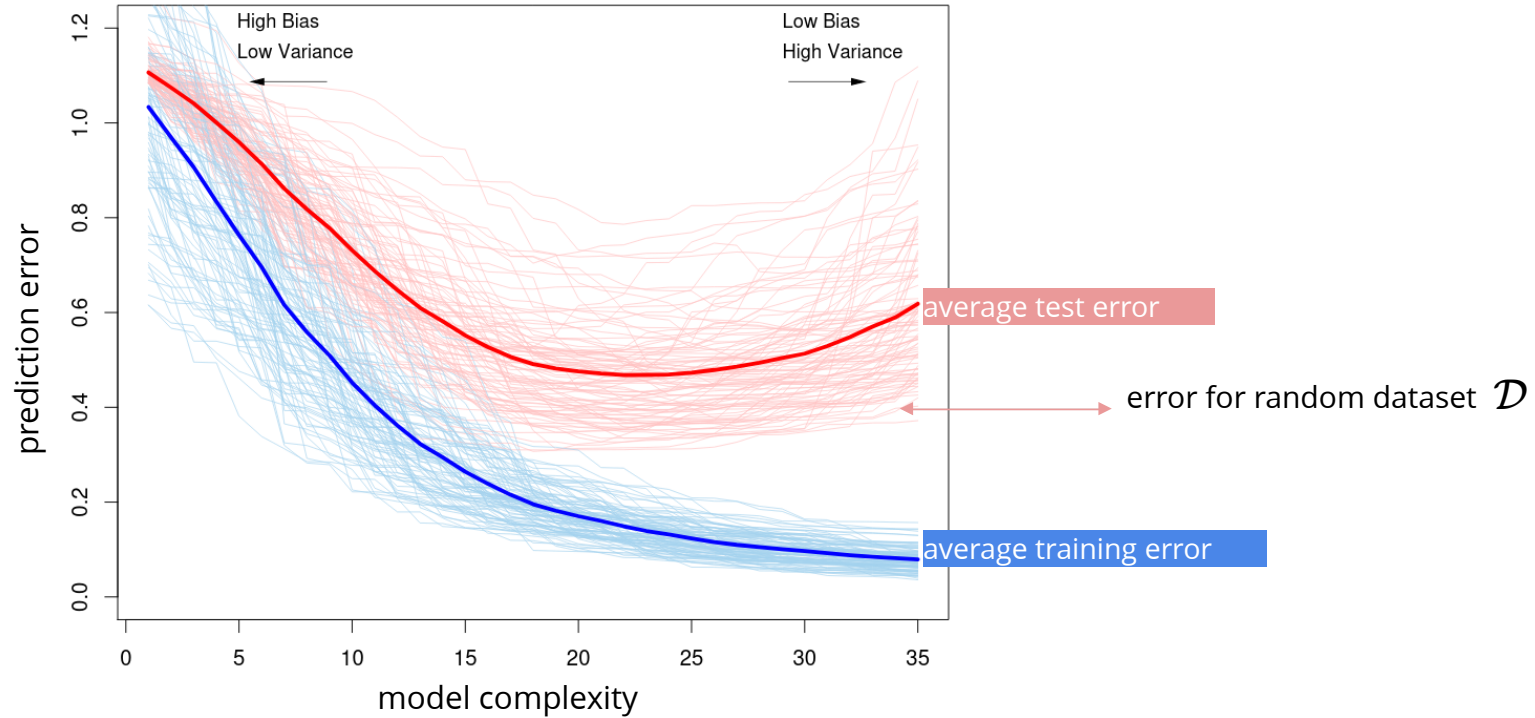
# Example: bias vs. variance

the lowest expected loss (test error) is somewhere between the two extremes



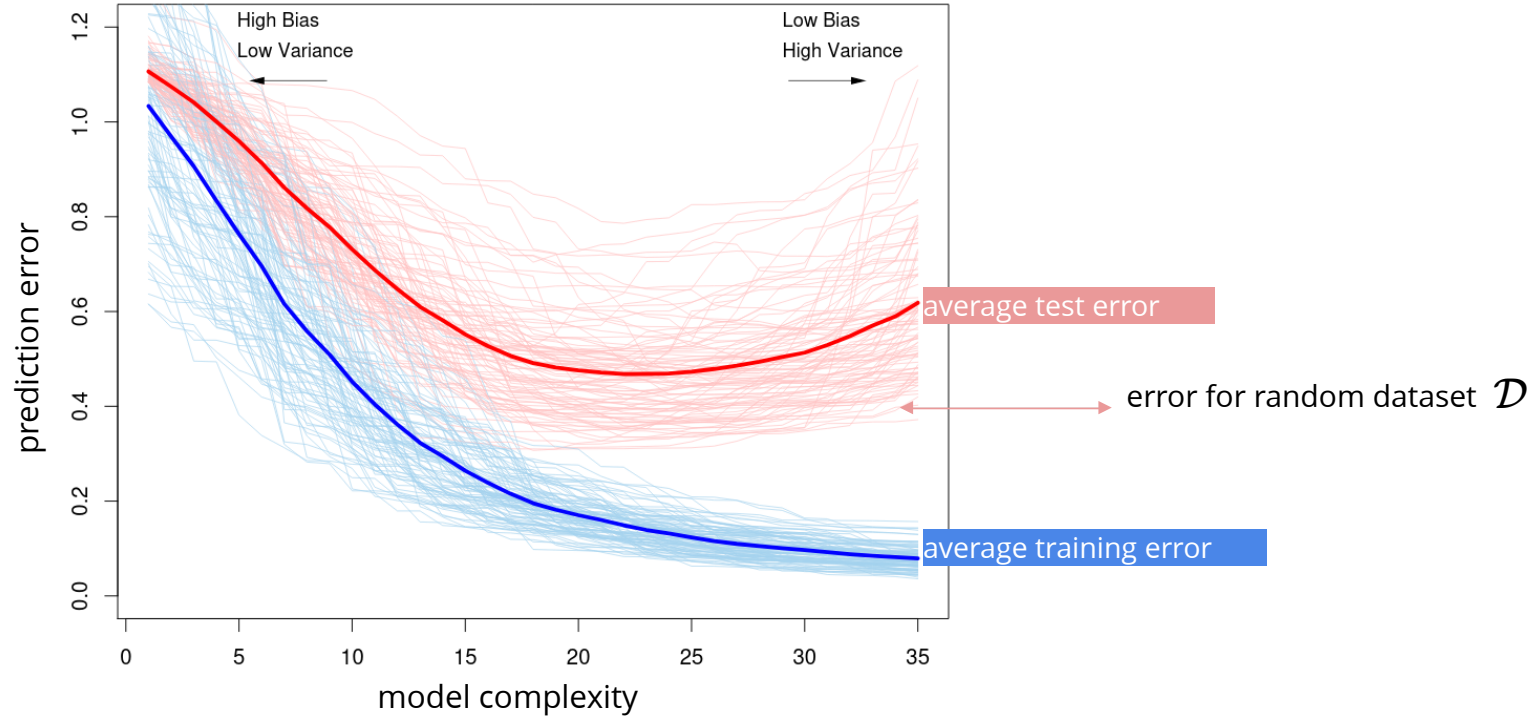
in reality, we don't have access to the true model  
how to decide which model to use?

# Big picture!



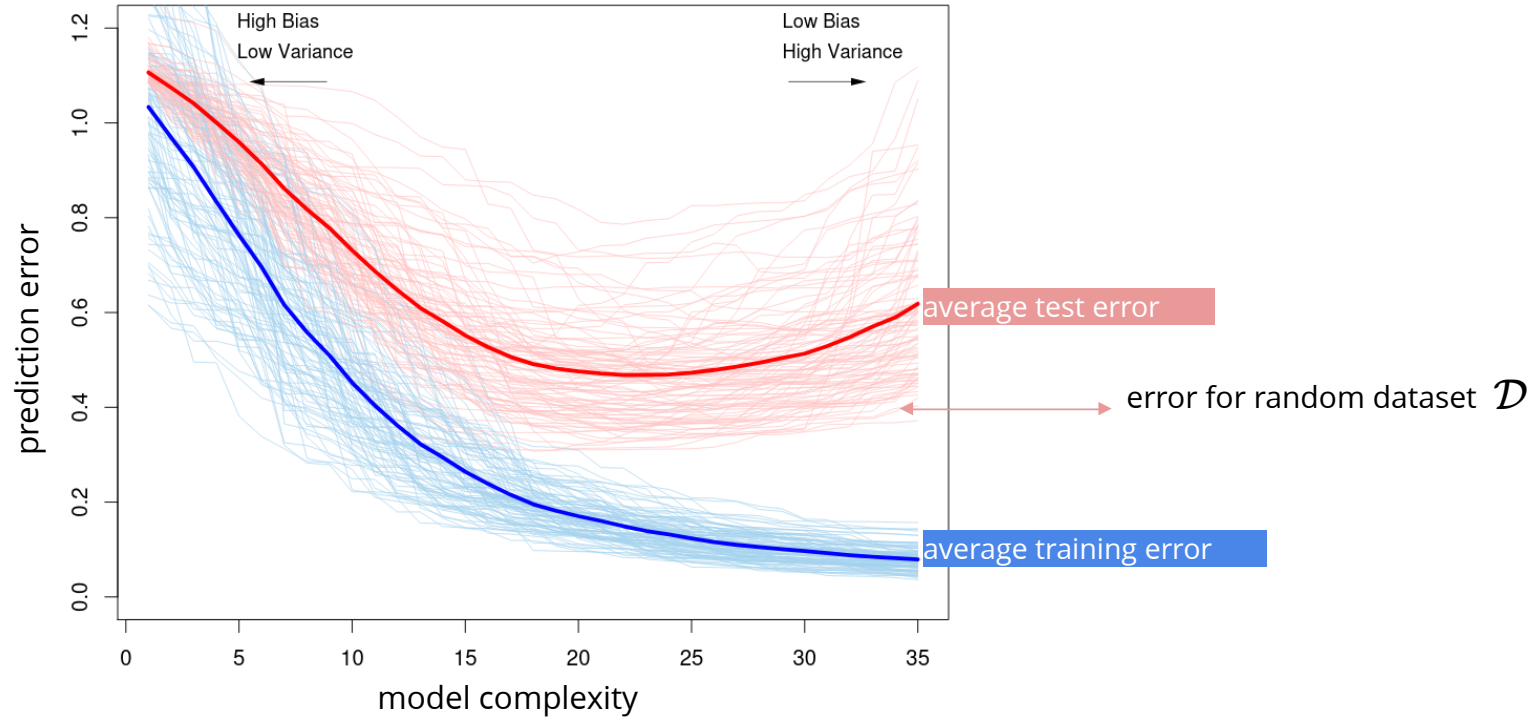
# Big picture!

high variance in more complex models means that test and training error can be very different



# Big picture!

high variance in more complex models means that **test and training error can be very different**  
high bias in simplistic models means that **training error can be high**

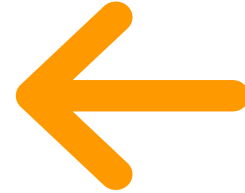




# Model selection

how to pick the model with lowest expected loss / test error?

use a **validation set** (and a separate test set for final assessment)



**regularization** bound the test error by bounding

- training error
- model complexity

# Model selection

how to pick the model with lowest expected loss / test error?

USE a **validation set** (and a separate test set for final assessment)



**regularization** bound the test error by bounding

- training error
- model complexity

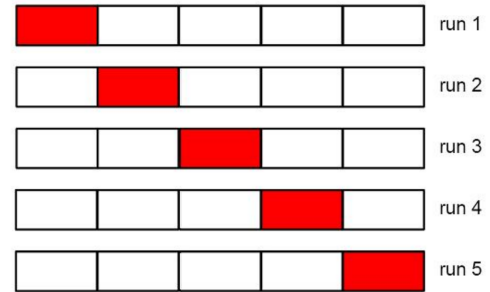
in the end we may have to use a validation set to find the right amount of regularization

# Cross validation

getting a more reliable estimate of test error using validation set

## K-fold cross validation(CV)

- randomly partition the data into  $K$  *folds*
- use  $K-1$  for training, and 1 for validation
- report average/std of the validation error over all folds

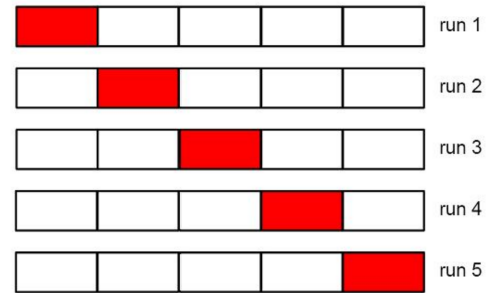


# Cross validation

getting a more reliable estimate of test error using validation set

## K-fold cross validation(CV)

- randomly partition the data into  $K$  *folds*
- use  $K-1$  for training, and 1 for validation
- report average/std of the validation error over all folds



leave-one-out CV: extreme case of  $k=N$

# Cross validation

getting a more reliable estimate of test error using validation set

## K-fold cross validation(CV)

- randomly partition the data into  $k$  *folds*
- use  $k-1$  for training, and 1 for validation
- report average/std of the validation error over all folds

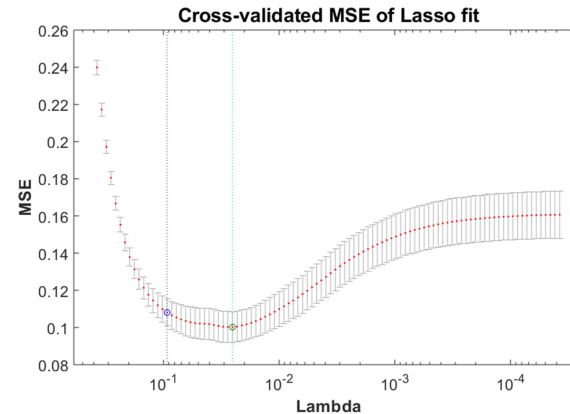


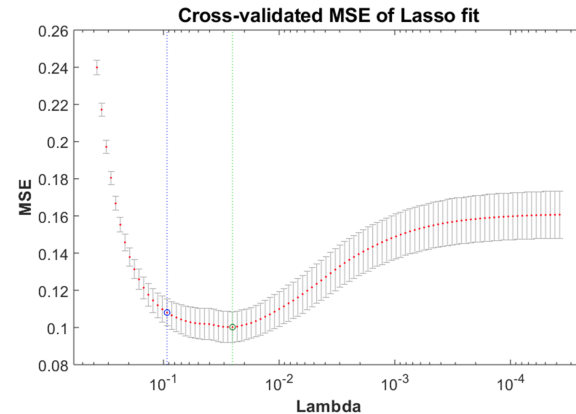
image credit: Thanh Nguyen et al'19

# Cross validation

getting a more reliable estimate of test error using validation set

## K-fold cross validation(CV)

- randomly partition the data into  $k$  *folds*
- use  $k-1$  for training, and 1 for validation
- report average/std of the validation error over all folds



once the hyper-parameters are selected, we can use the whole set for training  
use test set for the **final** assessment

image credit: Thanh Nguyen et al'19

# Evaluation

**evaluation metric** can be different from the optimization objective

**confusion matrix** is a CxC table that compares truth-vs-prediction

for **binary classification**:

	Truth		$\Sigma$
Result	TP	FP	RP
	FN	TN	RN
$\Sigma$	P	N	

# Evaluation

**evaluation metric** can be different from the optimization objective

**confusion matrix** is a CxC table that compares truth-vs-prediction

for **binary classification**:

	Truth		$\Sigma$
Result	TP	FP	RP
	FN	TN	RN
$\Sigma$	P	N	

some **evaluation metrics**  
(based on the confusion table)

$$Accuracy = \frac{TP+TN}{P+N}$$

$$Error\ rate = \frac{FP+FN}{P+N}$$

$$Precision = \frac{TP}{RP}$$

$$Recall = \frac{TP}{P}$$

$$F_1\ score = 2 \frac{Precision \times Recall}{Precision + Recall}$$



# Evaluation

**evaluation metric** can be different from the optimization objective  
**confusion matrix** is a CxC table that compares truth-vs-prediction

for **binary classification**:

Result	Truth		Σ
	TP	FP	RP
	FN	TN	RN
Σ	P	N	

some **evaluation metrics**  
(based on the confusion table)

$$Accuracy = \frac{TP+TN}{P+N}$$

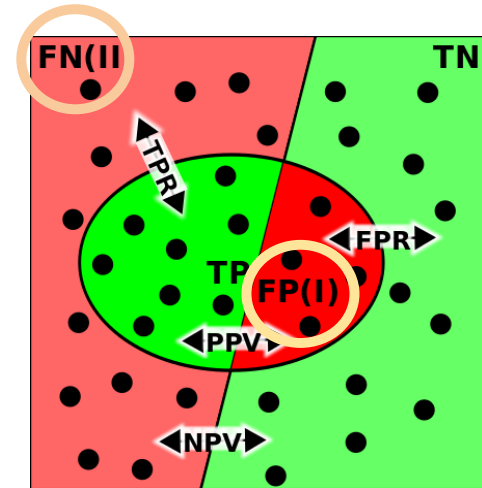
$$Error\ rate = \frac{FP+FN}{P+N}$$

$$Precision = \frac{TP}{RP}$$

$$Recall = \frac{TP}{P}$$

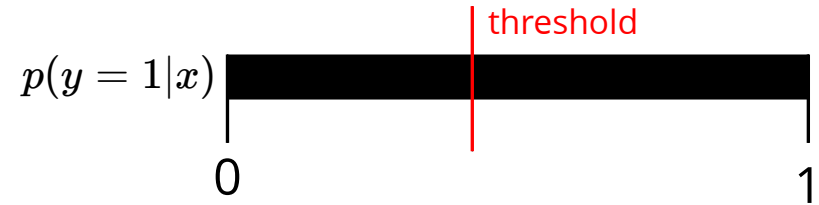
$$F_1\ score = 2 \frac{Precision \times Recall}{Precision + Recall}$$

type I vs type II error



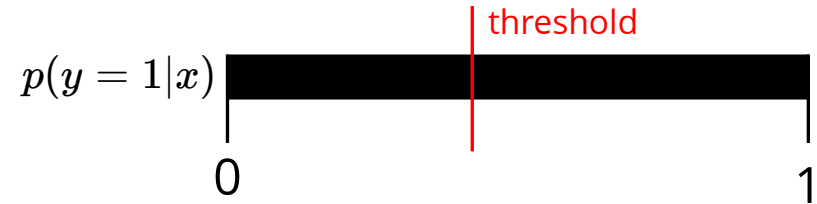
# Evaluation

if we produce class score (probability)  
we can trade-off between type I & type II error



# Evaluation

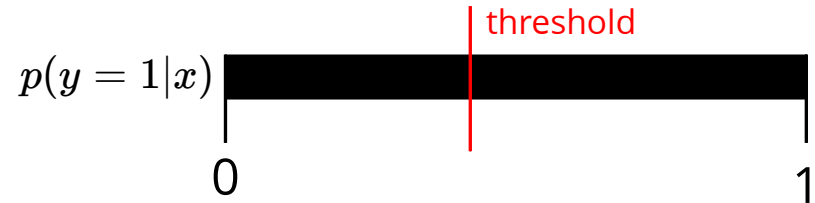
if we produce class score (probability)  
we can trade-off between type I & type II error



**goal:** evaluate class scores/probabilities (independent of choice of threshold)

# Evaluation

if we produce class score (probability)  
we can trade-off between type I & type II error

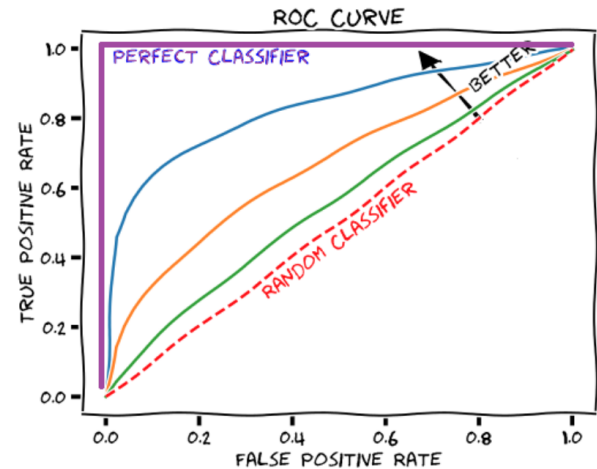


**goal:** evaluate class scores/probabilities (independent of choice of threshold)

Receiver Operating Characteristic **ROC curve**

**TPR** = TP/P (**recall**, sensitivity)

**FPR** = FP/N (**fallout**, false alarm)



# Summary

- complex models can have very different training and test error (*generalization gap*)
- regularization bounds this gap by penalizing model complexity
  - L1 & L2 regularization
  - probabilistic interpretation: different priors on weights
  - L1 produces sparse solutions (useful for feature selection)

# Summary

- complex models can have very different training and test error (*generalization gap*)
- regularization bounds this gap by penalizing model complexity
  - L1 & L2 regularization
  - probabilistic interpretation: different priors on weights
  - L1 produces sparse solutions (useful for feature selection)
- bias-variance trade off:
  - formalizes the relation between
    - training error (bias)
    - complexity (variance) and
    - and the test error (bias + variance)
  - not so elegant beyond L2 loss

# Summary

- complex models can have very different training and test error (*generalization gap*)
- regularization bounds this gap by penalizing model complexity
  - L1 & L2 regularization
  - probabilistic interpretation: different priors on weights
  - L1 produces sparse solutions (useful for feature selection)
- bias-variance trade off:
  - formalizes the relation between
    - training error (bias)
    - complexity (variance) and
    - and the test error (bias + variance)
  - not so elegant beyond L2 loss
- (cross) validation for model selection