

---

# COMP 551 – Applied Machine Learning

## Lecture 4: Linear classification

---

**Instructor:** Joelle Pineau ([jpineau@cs.mcgill.ca](mailto:jpineau@cs.mcgill.ca))

**Class web page:** [www.cs.mcgill.ca/~jpineau/comp551](http://www.cs.mcgill.ca/~jpineau/comp551)

Unless otherwise noted, all material posted for this course are copyright of the instructor, and cannot be reused or reposted without the instructor's written permission.

---

---

# Today's Quiz

---

1. What is meant by the term *overfitting*? What can cause overfitting? How can one avoid overfitting?
2. Which of the following increases the chances of overfitting (assuming everything else is held constant):
  - a) **Reducing** the size of the **training set**.
  - b) **Increasing** the size of the **training set**.
  - c) **Reducing** the size of the **test set**.
  - d) **Increasing** the size of the **test set**.
  - e) **Reducing** the number of **features**.
  - f) **Increasing** the number of **features**.

---

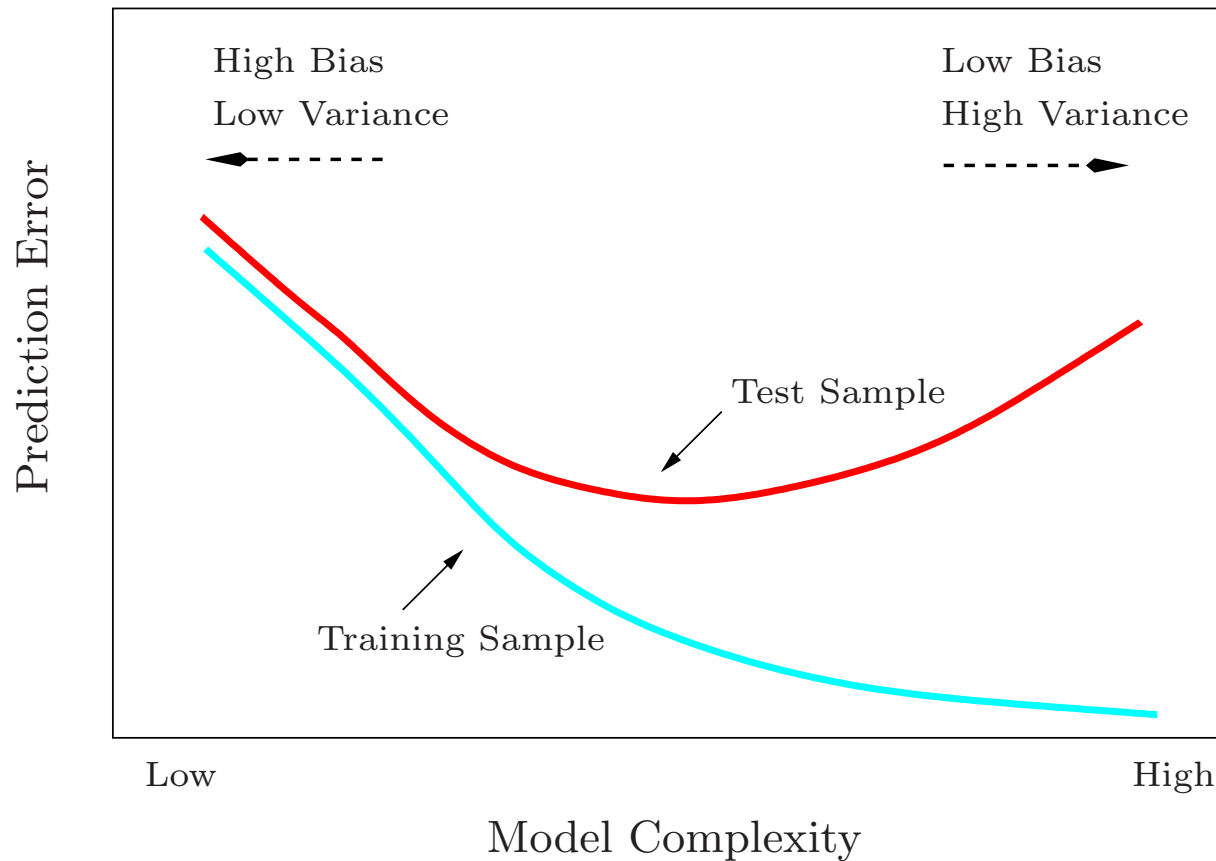
# Evaluation

---

- Use **cross-validation** for **model selection**.
- **Training set** is used to select a hypothesis  $f$  from a class of hypotheses  $F$  (e.g. regression of a given degree).
- **Validation set** is used to **compare best  $f$  from each hypothesis class** across different classes (e.g. different degree regression).
  - Must be untouched during the process of looking for  $f$  within a class  $F$ .
- **Test set**: Ideally, a separate set of (labeled) data is withheld to get a true estimate of the generalization error.  
(Often the “validation set” is called “test set”, without distinction.)

# Validation vs Train error

[From Hastie et al. textbook]



**FIGURE 2.11.** Test and training error as a function of model complexity.

---

# Bias vs Variance

---

- Gauss-Markov Theorem says:

The least-squares estimates of the parameters  $\mathbf{w}$  have the **smallest variance** among all linear **unbiased** estimates.

- **Insight:** Find **lower variance** solution, at the expense of **some bias**.

E.g. Include **penalty for model complexity** in error to reduce overfitting.

$$Err(\mathbf{w}) = \sum_{i=1:n} (y_i - \mathbf{w}^T \mathbf{x}_i)^2 + \lambda |model\_size|$$

$\lambda$  is a hyper-parameter that controls penalty size.

---

# Ridge regression (aka L2-regularization)

---

- Constrains the weights by imposing a penalty on their size:

$$\hat{\mathbf{w}}^{ridge} = \underset{\mathbf{w}}{\operatorname{argmin}} \left\{ \sum_{i=1:n} (y_i - \mathbf{w}^T \mathbf{x}_i)^2 + \lambda \sum_{j=0:m} w_j^2 \right\}$$

where  $\lambda$  can be selected manually, or by cross-validation.

- Do a little algebra to get the solution:  $\hat{\mathbf{w}}^{ridge} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{Y}$ 
  - The ridge solution is **not equivariant** under scaling of the data, so typically **need to normalize the inputs** first.
  - Ridge gives a smooth solution, effectively shrinking the weights, but drives few weights to 0.

---

# Lasso regression (aka L1-regularization)

---

- Constrains the weights by penalizing the absolute value of their size:

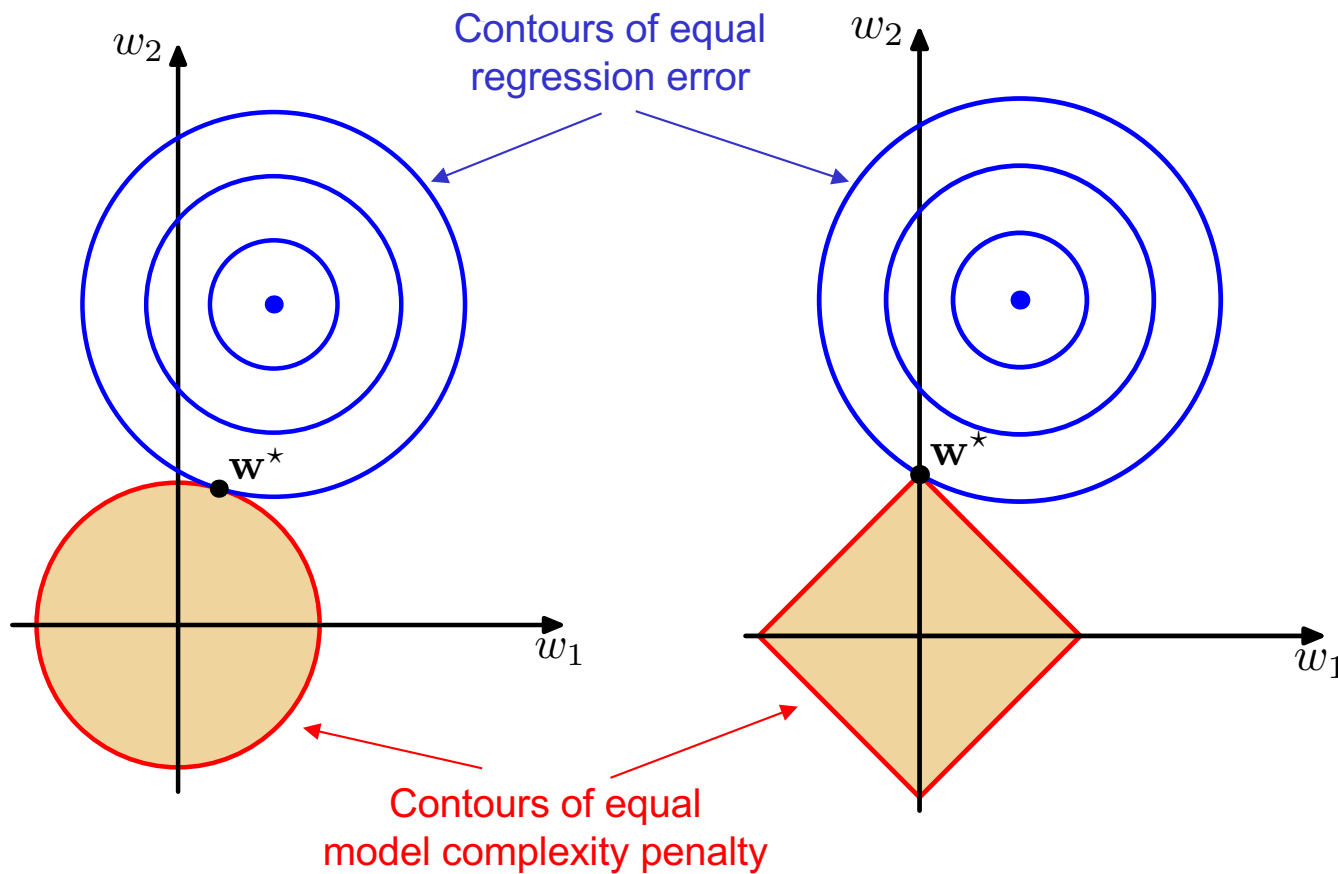
$$\hat{\mathbf{w}}^{\text{lasso}} = \underset{\mathbf{w}}{\operatorname{argmin}} \left\{ \sum_{i=1:n} (y_i - \mathbf{w}^T \mathbf{x}_i)^2 + \lambda \sum_{j=1:m} |w_j| \right\}$$

- Now the objective is non-linear in the output  $y$ , and there is no closed-form solution. **Need to solve a quadratic programming problem instead.**
  - More computationally expensive than Ridge regression.
  - Effectively sets the weights of less relevant input features to zero.

# Comparing Ridge and Lasso

## Ridge

## Lasso





---

# A quick look at evaluation functions

---

- We call  $L(Y, f_w(x))$  the loss function.

- Least-square / Mean squared-error (MSE) loss:

$$L(Y, f_w(X)) = \sum_{i=1:n} (y_i - \mathbf{w}^T x_i)^2$$

- Other loss functions?

- Absolute error loss:

$$L(Y, f_w(X)) = \sum_{i=1:n} |y_i - \mathbf{w}^T x_i|$$

- 0-1 loss (for classification):

$$L(Y, f_w(X)) = \sum_{i=1:n} I(y_i \neq f_w(x_i))$$

- Different loss functions make **different assumptions**.

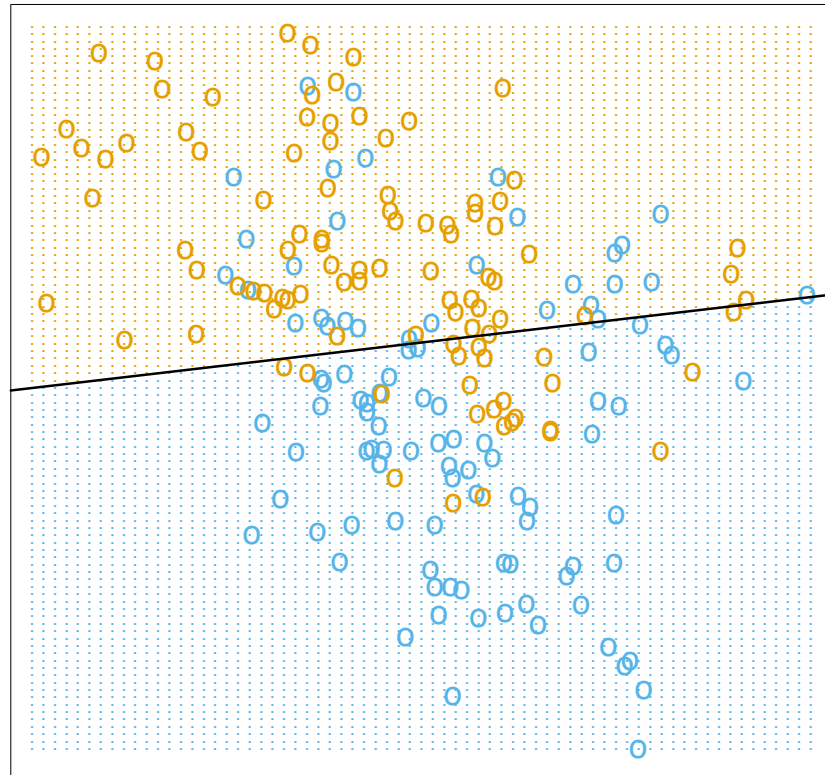
- Squared error loss assumes the data can be approximated by a global linear model with Gaussian noise.

---

# Next: Linear models for classification

---

Linear Regression of 0/1 Response



**FIGURE 2.1.** A classification example in two dimensions. The classes are coded as a binary variable (BLUE = 0, ORANGE = 1), and then fit by linear regression. The line is the decision boundary defined by  $x^T \hat{\beta} = 0.5$ . The orange shaded region denotes that part of input space classified as ORANGE, while the blue region is classified as BLUE.

# Classification problems

Given data set  $D = \langle x_i, y_i \rangle, i=1:n$ , with discrete  $y_i$ , find a hypothesis which “best fits” the data.

- If  $y_i \in \{0, 1\}$  this is **binary** classification.
- If  $y_i$  can take more than two values, the problem is called **multi-class** classification.

Fisher's Iris Data

Sepal length ↕	Sepal width ↕	Petal length ↕	Petal width ↕	Species ↕
5.1	3.5	1.4	0.2	<i>I. setosa</i>
4.9	3.0	1.4	0.2	<i>I. setosa</i>
4.7	3.2	1.3	0.2	<i>I. setosa</i>
4.6	3.1	1.5	0.2	<i>I. setosa</i>
5.0	3.6	1.4	0.2	<i>I. setosa</i>
5.4	3.9	1.7	0.4	<i>I. setosa</i>
4.6	3.4	1.4	0.3	<i>I. setosa</i>
5.0	3.4	1.5	0.2	<i>I. setosa</i>
4.4	2.9	1.4	0.2	<i>I. setosa</i>
4.9	3.1	1.5	0.1	<i>I. setosa</i>
5.4	3.7	1.5	0.2	<i>I. setosa</i>
4.8	3.4	1.6	0.2	<i>I. setosa</i>
4.8	3.0	1.4	0.1	<i>I. setosa</i>
4.3	3.0	1.1	0.1	<i>I. setosa</i>
5.8	4.0	1.2	0.2	<i>I. setosa</i>
5.7	4.4	1.5	0.4	<i>I. setosa</i>
5.4	3.9	1.3	0.4	<i>I. setosa</i>
5.1	3.5	1.4	0.3	<i>I. setosa</i>
5.7	3.8	1.7	0.3	<i>I. setosa</i>
5.1	3.8	1.5	0.3	<i>I. setosa</i>
5.4	3.4	1.7	0.2	<i>I. setosa</i>
5.1	3.7	1.5	0.4	<i>I. setosa</i>
4.6	3.6	1.0	0.2	<i>I. setosa</i>

---

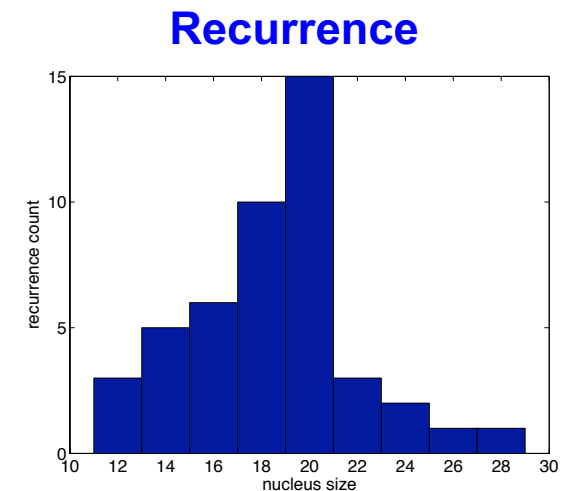
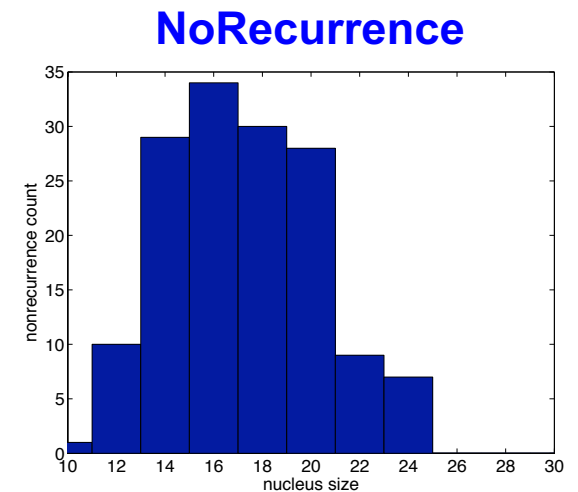
# Applications of classification

---

- Text classification (spam filtering, news filtering, building web directories, etc.)
- Image classification (face detection, object recognition, etc.)
- Prediction of cancer recurrence.
- Financial forecasting.
- Many, many more!

# Simple example

- Given “nucleus size”, predict cancer recurrence.
- Univariate input:  $X$  = nucleus size.
- Binary output:  $Y = \{\text{NoRecurrence} = 0; \text{Recurrence} = 1\}$
- Try: **Minimize the least-square error.**



# Predicting a class from linear regression

- Here red line is:  $Y' = X (X^T X)^{-1} X^T Y$

- How to get a binary output?

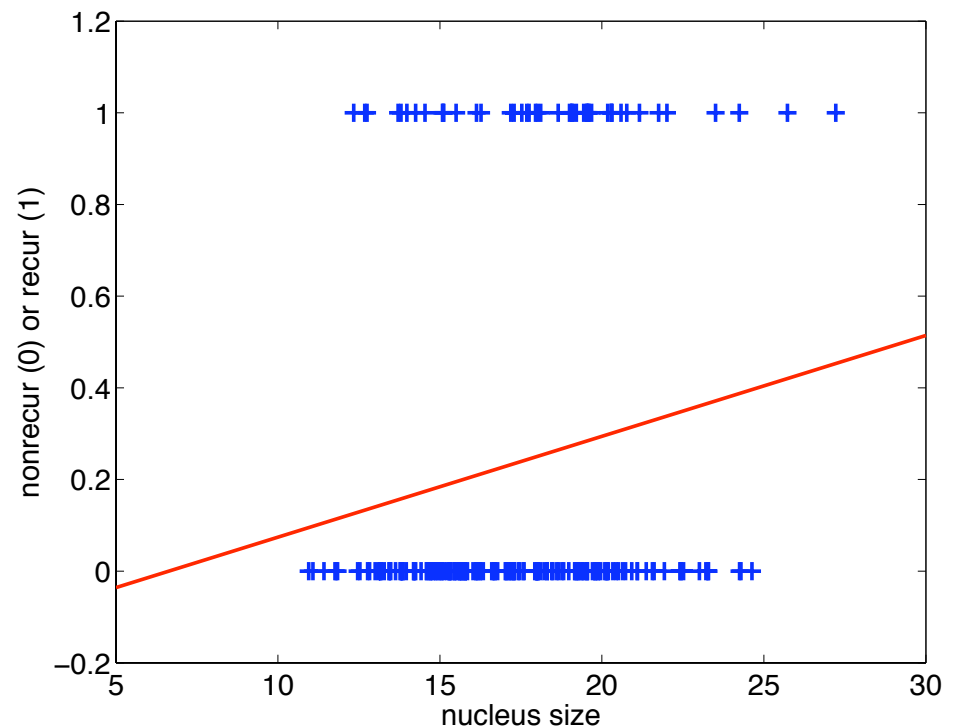
1. Threshold the output:

$\{ y \leq t$  for NoRecurrence,  
 $y > t$  for Recurrence}

2. Interpret output as probability:

$y = Pr(\text{Recurrence})$

- Can we find a better model?



---

# Modeling for binary classification

---

- Two probabilistic approaches:
  1. **Discriminative learning**: Directly estimate  $P(y|x)$ .
  2. **Generative learning**: Separately model  $P(x|y)$  and  $P(y)$ . Use Bayes rule, to estimate  $P(y|x)$ :

$$P(y = 1 | x) = \frac{P(x | y = 1)P(y = 1)}{P(x)}$$

---

# Probabilistic view of discriminative learning

---

- Suppose we have 2 classes:  $y \in \{0, 1\}$
- What is the probability of a given input  $x$  having class  $y = 1$ ?
- Consider Bayes rule:

$$\begin{aligned} P(y = 1 | x) &= \frac{P(x, y = 1)}{P(x)} = \frac{P(x | y = 1)P(y = 1)}{P(x | y = 1)P(y = 1) + P(x | y = 0)P(y = 0)} \\ &= \frac{1}{1 + \frac{P(x | y = 0)P(y = 0)}{P(x | y = 1)P(y = 1)}} = \frac{1}{1 + \exp(\ln \frac{P(x | y = 0)P(y = 0)}{P(x | y = 1)P(y = 1)})} = \frac{1}{1 + \exp(-a)} = \sigma \end{aligned}$$

where  $a = \ln \frac{P(x | y = 1)P(y = 1)}{P(x | y = 0)P(y = 0)} = \ln \frac{P(y = 1 | x)}{P(y = 0 | x)}$  (By Bayes rule;  $P(x)$  on top and bottom cancels out.)

- Here  $\sigma$  has a special form, called the **logistic function** and  $a$  is the log-odds ratio of data being class 1 vs. class 0.



---

# Discriminative learning: Logistic regression

---

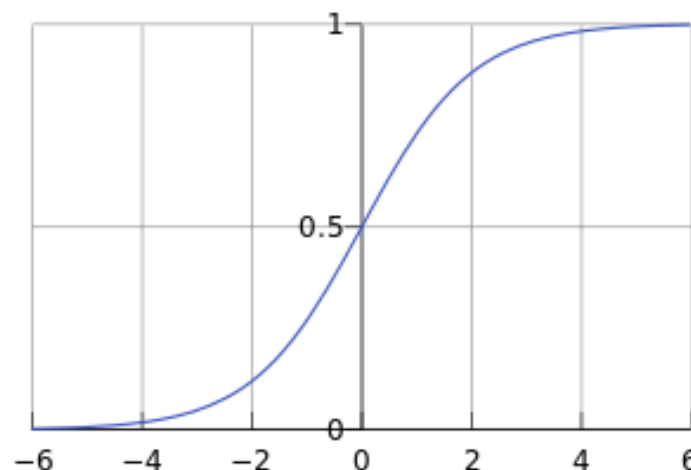
- Idea: Directly model the log-odds with a linear function:

$$a = \ln \frac{P(x | y = 1)P(y = 1)}{P(x | y = 0)P(y = 0)} = w_0 + w_1x_1 + \dots + w_mx_m$$

- The **decision boundary** is the set of points for which  $a=0$ .
- The **logistic** function (= sigmoid curve):  $\sigma(\mathbf{w}^T\mathbf{x}) = 1 / (1 + e^{-\mathbf{w}^T\mathbf{x}})$

How do we find the weights?

Need an optimization function.



---

# Fitting the weights

---

- Recall:  $\sigma(\mathbf{w}^T \mathbf{x}_i)$  is the probability that  $y_i=1$  (given  $\mathbf{x}_i$ )  
 $1-\sigma(\mathbf{w}^T \mathbf{x}_i)$  be the probability that  $y_i = 0$ .
- For  $y \in \{0, 1\}$ , the **likelihood** function,  $Pr(x_1, y_1, \dots, x_n, y_n | \mathbf{w})$ , is:  
$$\prod_{i=1:n} \sigma(\mathbf{w}^T \mathbf{x}_i)^{y_i} (1 - \sigma(\mathbf{w}^T \mathbf{x}_i))^{(1-y_i)}$$
 (samples are i.i.d.)
- **Goal:** Minimize the **log-likelihood** (also called *cross-entropy error function*):  
$$- \sum_{i=1:n} y_i \log(\sigma(\mathbf{w}^T \mathbf{x}_i)) + (1-y_i) \log(1-\sigma(\mathbf{w}^T \mathbf{x}_i))$$

---

# Gradient descent for logistic regression

---

- Error fn:  $Err(\mathbf{w}) = - [ \sum_{i=1:n} y_i \log(\sigma(\mathbf{w}^T \mathbf{x}_i)) + (1-y_i) \log(1-\sigma(\mathbf{w}^T \mathbf{x}_i)) ]$

- Take the derivative:

$$\frac{\partial Err(\mathbf{w})}{\partial \mathbf{w}} = - [ \sum_{i=1:n} y_i \log(\sigma(\mathbf{w}^T \mathbf{x}_i)) + (1-y_i) \log(1-\sigma(\mathbf{w}^T \mathbf{x}_i)) ]$$

$\downarrow \delta \log(\sigma) / \delta w = 1/\sigma$

$$\frac{\partial Err(\mathbf{w})}{\partial \mathbf{w}} = - [ \sum_{i=1:n} y_i (1/\sigma(\mathbf{w}^T \mathbf{x}_i)) (1-\sigma(\mathbf{w}^T \mathbf{x}_i)) \sigma(\mathbf{w}^T \mathbf{x}_i) \mathbf{x}_i + \dots ]$$

---

# Gradient descent for logistic regression

---

- Error fn:  $Err(\mathbf{w}) = - [ \sum_{i=1:n} y_i \log(\sigma(\mathbf{w}^T \mathbf{x}_i)) + (1-y_i) \log(1-\sigma(\mathbf{w}^T \mathbf{x}_i)) ]$

- Take the derivative:

$$\frac{\partial Err(\mathbf{w})}{\partial \mathbf{w}} = - [ \sum_{i=1:n} y_i \frac{1}{\sigma(\mathbf{w}^T \mathbf{x}_i)} (1-\sigma(\mathbf{w}^T \mathbf{x}_i)) \sigma(\mathbf{w}^T \mathbf{x}_i) \mathbf{x}_i + \dots ]$$

$$\delta\sigma/\delta\mathbf{w} = \sigma(1-\sigma)$$


---

# Gradient descent for logistic regression

---

- Error fn:  $Err(\mathbf{w}) = - [ \sum_{i=1:n} y_i \log(\sigma(\mathbf{w}^T \mathbf{x}_i)) + (1-y_i) \log(1-\sigma(\mathbf{w}^T \mathbf{x}_i)) ]$

- Take the derivative:

$$\frac{\partial Err(\mathbf{w})}{\partial \mathbf{w}} = - [ \sum_{i=1:n} y_i \frac{1}{\sigma(\mathbf{w}^T \mathbf{x}_i)} (1-\sigma(\mathbf{w}^T \mathbf{x}_i)) \sigma(\mathbf{w}^T \mathbf{x}_i) \mathbf{x}_i + \dots ]$$

$\delta w^T x / \delta w = x$



---

# Gradient descent for logistic regression

---

- Error fn:  $Err(\mathbf{w}) = - [ \sum_{i=1:n} y_i \log(\sigma(\mathbf{w}^T \mathbf{x}_i)) + (1-y_i) \log(1-\sigma(\mathbf{w}^T \mathbf{x}_i)) ]$

- Take the derivative:

$$\frac{\partial Err(\mathbf{w})}{\partial \mathbf{w}} = - [ \sum_{i=1:n} y_i (1/\sigma(\mathbf{w}^T \mathbf{x}_i))(1-\sigma(\mathbf{w}^T \mathbf{x}_i)) \sigma(\mathbf{w}^T \mathbf{x}_i) \mathbf{x}_i + (1-y_i)(1/(1-\sigma(\mathbf{w}^T \mathbf{x}_i)))(1-\sigma(\mathbf{w}^T \mathbf{x}_i))\sigma(\mathbf{w}^T \mathbf{x}_i)(-1) \mathbf{x}_i ]$$

$\delta(1-\sigma)/\delta w = (1-\sigma)\sigma(-1)$

---

# Gradient descent for logistic regression

---

- Error fn:  $Err(\mathbf{w}) = - [ \sum_{i=1:n} y_i \log(\sigma(\mathbf{w}^T \mathbf{x}_i)) + (1-y_i) \log(1-\sigma(\mathbf{w}^T \mathbf{x}_i)) ]$

- Take the derivative:

$$\begin{aligned} \partial Err(\mathbf{w}) / \partial \mathbf{w} &= - [ \sum_{i=1:n} y_i (1/\sigma(\mathbf{w}^T \mathbf{x}_i))(1-\sigma(\mathbf{w}^T \mathbf{x}_i)) \sigma(\mathbf{w}^T \mathbf{x}_i) \mathbf{x}_i + \\ &\quad (1-y_i)(1/(1-\sigma(\mathbf{w}^T \mathbf{x}_i)))(1-\sigma(\mathbf{w}^T \mathbf{x}_i))\sigma(\mathbf{w}^T \mathbf{x}_i)(-1) \mathbf{x}_i ] \\ &= - \sum_{i=1:n} \mathbf{x}_i (y_i (1-\sigma(\mathbf{w}^T \mathbf{x}_i)) - (1-y_i)\sigma(\mathbf{w}^T \mathbf{x}_i)) \\ &= - \sum_{i=1:n} \mathbf{x}_i (y_i - \sigma(\mathbf{w}^T \mathbf{x}_i)) \end{aligned}$$

---

# Gradient descent for logistic regression

---

- Error fn:  $Err(\mathbf{w}) = - [ \sum_{i=1:n} y_i \log(\sigma(\mathbf{w}^T \mathbf{x}_i)) + (1-y_i) \log(1-\sigma(\mathbf{w}^T \mathbf{x}_i)) ]$

- Take the derivative:

$$\begin{aligned} \partial Err(\mathbf{w}) / \partial \mathbf{w} &= - [ \sum_{i=1:n} y_i (1/\sigma(\mathbf{w}^T \mathbf{x}_i)) (1-\sigma(\mathbf{w}^T \mathbf{x}_i)) \sigma(\mathbf{w}^T \mathbf{x}_i) \mathbf{x}_i + \\ &\quad (1-y_i) (1/(1-\sigma(\mathbf{w}^T \mathbf{x}_i))) (1-\sigma(\mathbf{w}^T \mathbf{x}_i)) \sigma(\mathbf{w}^T \mathbf{x}_i) (-1) \mathbf{x}_i ] \\ &= - \sum_{i=1:n} \mathbf{x}_i (y_i (1-\sigma(\mathbf{w}^T \mathbf{x}_i)) - (1-y_i) \sigma(\mathbf{w}^T \mathbf{x}_i)) \\ &= - \sum_{i=1:n} \mathbf{x}_i (y_i - \sigma(\mathbf{w}^T \mathbf{x}_i)) \end{aligned}$$

- Now apply iteratively:

$$\mathbf{w}_{k+1} = \mathbf{w}_k + \alpha_k \sum_{i=1:n} \mathbf{x}_i (y_i - \sigma(\mathbf{w}_k^T \mathbf{x}_i))$$

- Can also apply other iterative methods, e.g. Newton's method, coordinate descent, L-BFGS, etc.



---

# Modeling for binary classification

---

- Two probabilistic approaches:
  1. **Discriminative learning**: Directly estimate  $P(y|x)$ .
  2. **Generative learning**: Separately model  $P(x|y)$  and  $P(y)$ . Use Bayes rule, to estimate  $P(y|x)$ :

$$P(y = 1 | x) = \frac{P(x | y = 1)P(y = 1)}{P(x)}$$

---

# What you should know

---

- Basic definition of linear classification problem.
- Derivation of logistic regression.
- Linear discriminant analysis: definition, decision boundary.
- Quadratic discriminant analysis: basic idea, decision boundary.
- LDA vs QDA pros/cons.
- Worth reading further:
  - Under some conditions, linear regression for classification and LDA are the same (Hastie et al., p.109-110).
  - Relation between Logistic regression and LDA (Hastie et al., 4.4.5)

---

# Final notes

---

You don't yet have a team for Project #1? => Use *myCourses*.

You don't yet have a plan for Project #1? => Start planning!

Feedback on tutorial 1?