Applied Machine Learning

Review

Reihaneh Rabbany



COMP 551 (Fall 2024)

1

ML Models



The core problem in ML is parameter estimation (aka model fitting), which requires solving an optimization problem of the loss/cost function

$$egin{aligned} J(w) &= rac{1}{N} \sum_{n=1}^N l(y^{(n)}, f(x^{(n)}; w)) \ & w^* &= rgmin_w J(w) \end{aligned}$$

Linear Regression: Logistic Regression: $\hat{y} = f_w(x) = w^ op x \ : \mathbb{R}^D o \mathbb{R}$ $\hat{y} = f_w(x) = \sigma(w^ op x): \mathbb{R}^D o \{0,1\}$ model: $J_w = rac{1}{N} \sum_n -y \log(\hat{y}^{(n)}) - (1-y^{(n)}) \log(1-\hat{y}^{(n)})$ $J_w = rac{1}{N} \sum_n rac{1}{2} (y^{(n)} - \hat{y}^{(n)})^2$ cost function: partial derivatives: $rac{\partial}{\partial w_d}J_w = rac{1}{N}\sum_n (\hat{y}^{(n)}-y^{(n)})x_d^{(n)}$ use SGD to find w^* $abla J(w) = rac{1}{N}\sum_n (\hat{y}^{(n)} - y^{(n)}) x^{(n)}$ gradient: vector of all partial derivatives: given $\nabla J(w)$ 2 Dx1

Linear model of regression

$$f_w(x) = w_0 + w_1 x_1 + \ldots + w_D x_D$$

model parameters or weights
bias or intercept
$$\frac{\text{simplification}}{\text{concatenate a 1 to } x \longrightarrow x = [1, x_1, \ldots, x_D]^{\top}$$
$$f_w(x) = w^{\top} x \qquad w = [w_0, w_1, \ldots, w_D]^{\top}$$

Logistic regression: model

recall the way we included a <code>bias</code> parameter $\,x=[1,x_1]\,$

the input feature is generated uniformly in [-5,5] for all the values less than 2 we have y=1 and y=0 otherwise

a good fit to this data is the one shown (green)

$$f_w(x) = \sigma(w^ op x) = rac{1}{1+e^{-w^ op}}$$

in the model shown w pprox [9.1, -4.5]

that is
$$\ \hat{y}=\sigma(-4.5x_1+9.1)$$

what is our model's decision boundary?



Example: binary classification



(petal width)

Gradient descent

an iterative algorithm for optimization

- starts from some $w^{\{0\}}$
- update using gradient $w^{\{t+1\}} \leftarrow w^{\{t\}} \alpha \nabla J(w^{\{t\}})$ steepest descent direction

converges to a local minima



cost function (for maximization : objective function)

t indicates the iteration step

learning rate

$$abla J(w) = [rac{\partial}{\partial w_1} J(w), \cdots rac{\partial}{\partial w_D} J(w)]^T$$

Minimum of a convex function

Convex functions are easier to minimize:

- critical points are global minimum
- gradient descent can find it

$$w^{\{t+1\}} \leftarrow w^{\{t\}} - lpha
abla J(w^{\{t\}})$$



Recognizing convex functions

the logistic regression cost function is convex in model parameters (w)

$$J(w) = rac{1}{N} \sum_{n=1}^{N} y^{(n)} \log \left(1 + e^{-w^{ op}x}
ight) + \left(1 - y^{(n)}
ight) \log \left(1 + e^{w^{ op}x}
ight)$$
same argument
checking second derivative
 $rac{\partial^2}{\partial z^2} \log(1 + e^z) = rac{e^{-z}}{(1 + e^{-z})^2} \ge 0$

sum of convex functions

Gradient for linear and logistic regression

in both cases: $abla J(w) = rac{1}{N}\sum_n x^{(n)}(\hat{y}^{(n)}-y^{(n)}) = rac{1}{N} X^ op (\hat{y}-y)$

linear regression: $\hat{y} = w^ op x$ logistic regression: $\hat{y} = \sigma(w^ op x)$

```
1 def gradient(x, y, w):
2 N,D = x.shape
3 yh = logistic(np.dot(x, w))
4 grad = np.dot(x.T, yh - y) / N
5 return grad
```

time complexity: $\mathcal{O}(ND)$

(two matrix multiplications)

compared to the direct solution for linear regression: $\mathcal{O}(ND^2 + D^3)$ gradient descent can be much faster for large $_D$

recall

example GD for linear regression

After 22 steps $w^{\{t+1\}} \leftarrow w^{\{t\}} - .01
abla J(w^{\{t\}})$



Learning rate α

Learning rate has a significant effect on GD



do a grid search usually between 0.001 to .1 to find the right value, look at the training curves

Stochastic Gradient Descent

we can write the cost function as an average over instances

 $J(w) = rac{1}{N} \sum_{n=1}^{N} J_n(w)$ cost for a single data-point *e.g. for linear regression*

$$J_n(w) = rac{1}{2} (w^T x^{(n)} - y^{(n)})^2$$

the same is true for the partial derivatives

$$rac{\partial}{\partial w_j}J(w) = rac{1}{N}\sum_{n=1}^N rac{\partial}{\partial w_j}J_n(w)$$

therefore $abla J(w) = \mathbb{E}_{\mathcal{D}}[
abla J_n(w)]$

Stochastic Gradient Descent

Idea: use stochastic approximations $\nabla J_n(w)$ in gradient descent

stochastic gradient update

 $w \leftarrow w - \alpha
abla J_{n}(w)$

the steps are "on average" in the right direction



each step is using gradient of a different cost, $J_n(w)$

each update is (1/N) of the cost of batch gradient

e.g., for linear regression $\mathcal{O}(D)$ $abla J_n(w) = x^{(n)}(w^ op x^{(n)} - y^{(n)})$ batch gradient update

 $w \leftarrow w - lpha
abla J(w)$

with small learning rate: guaranteed improvement at each step



 w_0

contour plot of the cost function

SGD for logistic regression example

logistic regression for Iris dataset (D=2 , $\, lpha = .1$)



stochastic gradient



Convergence of SGD

stochastic gradients are not zero even at the optimum w how to guarantee convergence?

idea: schedule to have a smaller learning rate over time

Robbins Monro

the sequence we use should satisfy: $\sum_{t=0}^{\infty} \alpha^{\{t\}} = \infty$ & otherwise for large $||w^{\{0\}} - w^*||$ we can't reach the minimum the steps should go to zero $\sum_{t=0}^{\infty} (\alpha^{\{t\}})^2 < \infty$

example
$$lpha^{\{t\}}=rac{10}{t}, lpha^{\{t\}}=t^{-.51}$$
 .





Minibatch SGD

use a minibatch to produce gradient estimates



Momentum

to help with oscillations:

- use a **running average** of gradients
- more recent gradients should have higher weights

$$\Delta w^{\{t\}} \leftarrow eta \Delta w^{\{t-1\}} + (1-eta)
abla J_{\mathbb{B}}(w^{\{t-1\}})$$

 $w^{\{t\}} \leftarrow w^{\{t-1\}} - \frac{\alpha \Delta w^{\{t\}}}{\alpha}$

momentum of 0 reduces to SGD common value > .9

is effectively an exponential moving average

$$\Delta w^{\{T\}} = \sum_{t=1}^{T} \beta^{T-t} (1-\beta) \nabla J_{\mathbb{B}}(w^{\{t\}})$$

there are other variations of momentum with similar idea



weight for the **most recent** gradient $(1 - \beta)$

Momentum



see the beautiful demo at Distill https://distill.pub/2017/momentum/

RMSprop (Root Mean Squared propagation)

solve the problem of diminishing step-size with Adagrad

• use exponential moving average instead of sum (similar to momentum)

instead of Adagrad: $S_d^{\{t\}} \leftarrow S_d^{\{t-1\}} + rac{\partial}{\partial w_d} J(w^{\{t-1\}})^2$

$$egin{aligned} S^{\{t\}} &\leftarrow \gamma S^{\{t-1\}} + (1-\gamma)
abla J(w^{\{t-1\}})^2 \ w^{\{t\}} &\leftarrow w_{\{t-1\}} - rac{lpha}{\sqrt{S^{\{t\}} + \epsilon}}
abla J(w^{\{t-1\}}) \end{aligned}$$
 ic

identical to Adagrad

note that $S^{\{t\}}$ here is a vector and with the square root is element-wise

Adam (Adaptive Moment Estimation)

two ideas so far:

- 1. use momentum to smooth out the oscillations
- 2. adaptive per-parameter learning rate

Adam combines the two:

both use exponential moving averages

 $\left| \begin{array}{l} M^{\{t\}} \leftarrow \beta_1 M^{\{t-1\}} + (1-\beta_1) \nabla J(w^{\{t-1\}}) & \text{identical to method of momentum} \\ S^{\{t\}} \leftarrow \beta_2 S^{\{t-1\}} + (1-\beta_2) \nabla J(w^{\{t-1\}})^2 & \text{identical to RMSProp} \\ w^{\{t\}} \leftarrow w^{\{t-1\}} - \frac{\alpha}{\sqrt{\hat{S}^{\{t\}}} + \epsilon} \hat{M}^{\{t\}} \end{array} \right|$

since M and S are initialized to be zero, at early stages they are biased towards zero

$$\hat{M}^{\{t\}} \gets rac{M^{\{t\}}}{1-eta_1^t} \qquad \hat{S}^{\{t\}} \gets rac{S^{\{t\}}}{1-eta_2^t}$$

for large time-steps it has no effect for small t, it scales up numerator

Loss, cost and generalization

assume we have a **model** $f:x\mapsto y$ for example $f:|{f 3}|{f \mapsto 3}|$

and we have a **loss function** that measures the error in our prediction $~\ell:y,\hat{y} o\mathbb{R}$

for example
$$igg| egin{array}{cc} \ell(y,\hat{y})=(y-\hat{y})^2 & ext{ for regression} \ \ell(y,\hat{y})=\mathbb{I}(y
eq\hat{y}) & ext{ for classification} \end{array}$$

we train our models to minimize **the cost function**:

$$J = \frac{1}{|\mathcal{D}_{\text{train}}|} \sum_{x,y \in \mathcal{D}_{\text{train}}} \ell(y, f(x))$$

We can drop this, why?
What we really care about is the **generalization error**: $\mathbb{E}_{x,y \sim p} \ \ell(y, f(x))$.

we can set aside part of the given data and use it to estimate generalization error

Bias-variance decomposition: Setup

decompose the generalization error to see the effect of bias and variance (for L2 loss) assume a true distribution $\,p(x,y)\,$

best prediction given L2 loss $\ \ f(x) = \mathbb{E}_p[y|x]$

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 generalization error (aka expected loss, expected risk)

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

all blue items are random variables

Bias-variance decomposition

what we care about is the generalization error

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

 $\begin{vmatrix} & | \\ & f(x) + \epsilon \\ \hat{f}_{\mathcal{D}}(x) + \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] \text{ add and subtract a term} \end{vmatrix}$

above simplifies to the following (the remaining terms are going to be zero)

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

variance bias^2 unavoidable noise error



Example: bias vs. variance



Validation set

what we really care about is the **generalization error**: $\mathbb{E}_{x,y\sim p}\;\ell(y,f(x))$

we can set aside part of the training data and use it to **estimate** the generalization error

train	training validation unseen (test)		
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0	0 0	00000000000000000000000000000000000000	

pick a hyper-parameter that gives us the best validation error

at the very end, we report the error on **test set**

validation and test error could be different because they use limited amount of data



how to estimate this?

0

Cross validation

- divide the (training + validation) data into L parts
- use one part for validation and L-1 for training



train

• report the test error for the final model

this is called **L-fold** cross-validation

in **leave-one-out** cross-validation L=N (only one instance is used for validation)



 $ar{e} = rac{1}{5} \sum_{i=1}^5 e_i$

 use the average validation error and its variance (uncertainty) to pick the best model

 e_{t}

27

Performance metrics for classification

Not all errors are the same

In particular in classification, we have different types of mistakes

false positive (type I) and false negative (type II)

example:

patient does not have disease but received positive diagnostic (Type I error) patient has disease but it was not detected (Type II error)

a message that is not spam is assigned to the spam folder (Type I error) a message that is spam appears in the regular folder (Type II error)

Performance metrics for classification

confusion matrix	Trut	h	Σ
Regult	TP	FP	RP
Itesuit	FN	ΤN	RN
Σ	Р	Ν	



 $Accuracy = \frac{TP+TN}{P+N}$ $Precision = \frac{TP}{RP}$ $Recall = \frac{TP}{P}$ sensitivity $F_1 score = 2 \frac{Precision \times Recall}{Precision + Recall}$ {Harmonic mean} $F_eta score = (1+eta^2) rac{Precision imes Recall}{eta^2 Precision + Recall}$ recall is β times more important compared to precision $Miss\ rate = rac{FN}{P}$ $Fallout = rac{FP}{N}$ $\begin{array}{l} \mathsf{Fallout} = \frac{FP}{N} \\ \mathsf{False} \ discovery \ r \\ Selectivity = \frac{TN}{N} \\ \mathsf{False} \ omission \ r \end{array}$ false positive rate False discovery rate = $\frac{FP}{RP}$ specificity $False \ omission \ rate = rac{FN}{RN} \ Negative \ predictive \ value = rac{TN}{RN}$

Trade-off between precision and recall

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

Receiver Operating Characteristic ROC curve, a function of threshold t

TPR(t) = TP(t)/P (**recall**, sensitivity at t, hit rate) **FPR(t)** = FP(t)/N (**fallout**, false alarm at t, type I error rate)

Area Under the Curve (**AUC**) is used as a threshold independent measure of quality of the classifier

 $AUC = \sum_t TPR(t)(FPR(t) - FPR(t-1))$, box-rule approximation



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 $\ (\Phi^ op \Phi) w^* = \Phi^ op y$ replacing X with Φ a (nonlinear) feature $\Phi = egin{bmatrix} \phi_1(x^{(1)}), & \phi_2(x^{(1)}), & \cdots, & \phi_D(x^{(1)}) \ \phi_1(x^{(2)}), & \phi_2(x^{(2)}), & \cdots, & \phi_D(x^{(2)}) \ dots & dots & \ddots & dots \ \phi_1(x^{(N)}), & \phi_2(x^{(N)}), & \cdots, & \phi_D(x^{(N)}) \end{bmatrix}$ one instance

Example: Gaussian bases





Example: Gaussian bases





Overfitting

which one of these models performs better at test time?



An observation

when overfitting, we sometimes see large weights





idea: penalize large parameter values

Ridge regression

also known as

L2 regularized linear least squares regression:

$$egin{aligned} J(w) &= rac{1}{2}||Xw-y||_2^2 + rac{\lambda}{2}||w||_2^2 \ &ert &ert &ert \ &ert &ert \ &er$$

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

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

 λ is a hyper-parameter (use a validation set or cross-validation to pick the best value)

Ridge regression

set the derivative to zero $J(w) = \frac{1}{2} \sum_{x,y \in D} (y - w^{\top}x)^2 + \frac{\lambda}{2} w^{\top}w$ $\nabla J(w) = \sum_{x,y \in D} x(w^{\top}x - y) + \lambda w$ $= X^{\top}(Xw - y) + \lambda w = 0$

linear system of equations $(X^ op X + \lambda \mathbf{I})w = X^ op y$

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

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

the only part different due to regularization

 λI makes it invertible, adds a small value to the diagonals $X^ op X$ we can have linearly dependent features the solution will be unique!

Probabilistic interpretation

recall linear regression & logistic regression maximize log-likelihood

 $w^{MLE} = rgmax_w \, p(y|X,w)$

linear regression
$$w^{MLE} = rg\max_w \prod_{x,y\in\mathcal{D}} \mathcal{N}ig(y|w^ op x,\sigma^2ig)$$

logistic regression $w^{MLE} = rg\max_w \prod_{x,y\in\mathcal{D}} \operatorname{Bernoulli}ig(y|\sigma(w^ op x))$

can we do Bayesian inference instead of maximum likelihood?

$$p(w|y,X) \propto p(w)p(y|w,X)$$
 posterior prior likelihood

Gaussian Prior

MAP estimate
$$w^{MAP} = rg\max_w \log p(y|X,w) + \frac{\log p(w)}{\operatorname{prior}}$$

assume independent zero-mean Gaussians

$$\log p(w) = \log \prod_{d=1}^D \mathcal{N}(w_d|0, au^2) = -\sum_d rac{w^2}{2 au^2} + ext{const.}$$

does not depend on w so it doesn't affect the optimization

lets call $\; rac{1}{ au^2} o \lambda$ then we get the L2 regularization penalty $rac{\lambda}{2}||w||_2^2$

smaller variance of the prior au gives larger regularization λ



 $\mathcal{N}(\mu,\sigma) = rac{1}{\sigma\sqrt{2\pi}}e^{-rac{1}{2}(rac{x-\mu}{\sigma})}$

Probabilistic interpretation for Linear Regression

idea given the dataset
$$\mathcal{D} = \{(x^{(1)}, y^{(1)}), \dots, (x^{(N)}, y^{(N)})\}$$

learn a probabilistic model p(y|x;w)



consider
$$p(y|x;w)$$
 with the following form $p_w(y \mid x) = \mathcal{N}(y \mid w^\top x, \sigma^2) = rac{1}{\sqrt{2\pi\sigma^2}} e^{-rac{(y-w^\top x)^2}{2\sigma^2}}$ assume a fixed variance, say $\sigma^2 = 1$

Q: how to fit the model?

A: maximize the conditional likelihood!

Maximum likelihood & linear regression

cond. probability
$$p(y \mid x; w) = \mathcal{N}(y \mid w^{\top}x, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-w^{\top}x)^2}{2\sigma^2}}$$

 y
 y
 $w^T x$ likelihood $L(w) = \prod_{n=1}^{N} p(y^{(n)} \mid x^{(n)}; w)$
 \log
 \log likelihood $\ell(w) = \sum_n -\frac{1}{2\sigma^2} (y^{(n)} - w^{\top}x^{(n)})^2 + \text{constants}$
 \max -likelihood params. $w^* = \arg \max_w \ell(w) = \arg \min_w \frac{1}{2} \sum_n (y^{(n)} - w^{\top}x^{(n)})^2$
linear least squares!
whenever we use square loss, we are assuming Gaussian noise!

Laplace prior

another notable choice of prior is the Laplace distribution

minimizing negative log-likelihood $\longrightarrow \sum_{d} \log p(w_d) = -\sum_{d} \frac{1}{\beta} |w_d| = -\frac{1}{\beta} ||w||_1$

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

(least absolute shrinkage and selection operator)

L1 norm of w

42



$L_1 \text{ vs } L_2$ regularization

regularization path shows how $\{w_d\}$ change as we change λ Lasso produces sparse weights (many are zero, rather than small)



Probabilistic view for Logistic Regression

Interpret the prediction as class probability $\ \hat{y} = p_w(y = 1 \mid x) = \sigma(w^ op x)$

the log-ratio of class probabilities is linear

$$\log rac{\hat{y}}{1-\hat{y}} = \log rac{\sigma(w^ op x)}{1-\sigma(w^ op x)} = \log rac{1}{e^{-w^ op x}} = w^ op x$$

logit function is the inverse of logistic, read more here

We have a Bernoulli likelihood

$$p(y^{(n)} \mid x^{(n)};w) = ext{Bernoulli}(y^{(n)};\sigma(w^ op x^{(n)})) = \hat{y}^{(n)}{}^{y^{(n)}}(1-\hat{y}^{(n)})^{1-y^{(n)}}$$

conditional likelihood of the labels given the inputs $L(w) = \prod_{n=1}^{N} p(y^{(n)} \mid x^{(n)}; w) = \prod_{n=1}^{N} \hat{y}^{(n)} (1 - \hat{y}^{(n)})^{1-y^{(n)}}$

Maximum likelihood & logistic regression

likelihood
$$L(w) = \prod_{n=1}^N p(y^{(n)} \mid x^{(n)};w) = \prod_{n=1}^N \hat{y}^{(n)} {y^{(n)}} (1-\hat{y}^{(n)})^{1-y^{(n)}}$$

find w that maximizes

$$egin{aligned} & w^* = \max_w \sum_{n=1}^N \log p_w(y^{(n)} \mid x^{(n)};w) \ & = \max_w \sum_{n=1}^N y^{(n)} \log(\hat{y}^{(n)}) + (1-y^{(n)}) \log(1-\hat{y}^{(n)}) \end{aligned}$$

 $= \min_{w} J(w)$ the cross entropy cost function!

so using cross-entropy loss in logistic regression is maximizing conditional likelihood

we saw a similar interpretation for linear regression (L2 loss maximizes the conditional Gaussian likelihood)

Multiclass classification

using this probabilistic view we extend logistic regression to multiclass setting

binary classification: Bernoulli likelihood:

 $\begin{array}{ll} \operatorname{Bernoulli}(y \mid \hat{y}) = \hat{y}^{y}(1 - \hat{y})^{1 - y} & \text{subject to} & \hat{y} \in [0, 1] & \begin{cases} \hat{y} & y = 1 \\ 1 - \hat{y} & y = 0 \end{cases} \\ & \checkmark & \\ & \mathsf{using logistic function to ensure this } & \hat{y} = \sigma(z) = \sigma(w^{T}x) \end{cases}$

C classes: categorical likelihood

$$ext{Categorical}(y \mid \hat{y}) = \prod_{c=1}^C \hat{y}_c^{\mathbb{I}(y=c)} \quad ext{ subject to } \quad \sum_c \hat{y}_c = 1$$

achieved using softmax function

 $egin{cases} \hat{y}_1 & y=1 \ \hat{y}_2 & y=2 \ \dots & \hat{y}_2 & u=C \end{cases}$

Softmax

generalization of logistic to > 2 classes:

- **logistic**: $\sigma : \mathbb{R} \to (0, 1)$ produces a single probability
 - probability of the second class is $(1 \sigma(z))$
- softmax: $\mathbb{R}^C o \Delta_C$ recall: probability simplex $p \in \Delta_c o \sum_{c=1}^C p_c = 1$

$$\hat{y}_c = ext{softmax}(z)_c = rac{e^{zc}}{\sum_{c'=1}^C e^{z_{c'}}}$$
 so $\sum_c \hat{y} = 1$

example softmax
$$([1, 1, 2, 0]) = [\frac{e}{2e+e^2+1}, \frac{e}{2e+e^2+1}, \frac{e^2}{2e+e^2+1}, \frac{1}{2e+e^2+1}]$$

$$\operatorname{softmax}([10,100,-1])\approx [0,1,0]$$

if input values are large, softmax becomes similar to argmax similar to logistic this is also a squashing function

Multiclass classification

C classes: categorical likelihood

 $ext{Categorical}(y \mid \hat{y}) = \prod_{c=1}^C \hat{y}_c^{\mathbb{I}(y=c)}$ using softmax to enforce sum-to-one constraint

T

$$\hat{y}_c = ext{softmax}([w_1^{ op}x,\ldots,w_C^{ op}x])_c = rac{e^{w_c^{ op}x}}{\sum_{c'}e^{w_{c'}^{ op}x}}$$
 so we have on parameter **vector** for each class $w_1 = [w_{1,1},w_{1,2},\ldots w_{1,D}]$

to simplify equations we write $z_c = w_c^{+} x_c^{+}$

$$\hat{y}_c = ext{softmax}([z_1,\ldots,z_C])_c = rac{e^{z_c}}{\sum_{c'} e^{z_{c'}}}$$

Likelihood for multiclass classification

C classes: categorical likelihood

 $ext{Categorical}(y \mid \hat{y}) = \prod_{c=1}^C \hat{y}_c^{\mathbb{I}(y=c)}$ using softmax to enforce sum-to-one constraint

$$\hat{y}_c = ext{softmax}([z_1,\ldots,z_C])_c = rac{e^{z_c}}{\sum_{c'}e^{z_{c'}}}$$
 where $z_c = w_c^{ op}x$

substituting softmax in Categorical likelihood:

$$\begin{aligned} \mathsf{likelihood} \quad & L(\{w_c\}) = \prod_{n=1}^{N} \prod_{c=1}^{C} \mathsf{softmax}([z_1^{(n)}, \dots, z_C^{(n)}])_c^{\mathbb{I}(y^{(n)} = c)} \\ & = \prod_{n=1}^{N} \prod_{c=1}^{C} \left(\frac{e^{z_c^{(n)}}}{\sum_{c'} e^{z_c'}}\right)^{\mathbb{I}(y^{(n)} = c)} \end{aligned}$$

One-hot encoding

$$\begin{array}{ll} \text{likelihood} & L(\{w_c\}) = \prod_{n=1}^{N} \prod_{c=1}^{C} \left(\frac{e^{z_c^{(n)}}}{\sum_{c'} e^{z_{c'}^{(n)}}} \right)^{\mathbb{I}(y^{(n)}=c)} \\ \\ \text{log-likelihood} & \ell(\{w_c\}) = \sum_{n=1}^{N} \sum_{c=1}^{C} \mathbb{I}(y^{(n)}=c) (z_c^{(n)} - \log \sum_{c'} e^{z_c^{(n)}}) \end{array}$$

one-hot encoding for labels $y^{(n)} o [\mathbb{I}(y^{(n)}=1),\ldots,\mathbb{I}(y^{(n)}=C)]$ convert that feature into C binary features

Example: $y^{(n)} \in \{1,2,3\} \Rightarrow y^{(n)} \in \{[1,0,0],[0,1,0],[0,0,1]\}$

side note

we can also use this encoding for categorical **features**

$$x_d^{(n)}
ightarrow [\mathbb{I}(x_d^{(n)}=1),\ldots,\mathbb{I}(x_d^{(n)}=C)]$$

Food Na	me	Categ	orical #	Ca	lories	
Apple		1		95	95	
Chicken		2		23	231	
Broccoli		3		50	50	
Apple	Chi	icken	Brocco	oli	Calories	
1	0		0		95	
0	1		0		231	
0	0		1		50	

One-hot encoding

$$\begin{array}{ll} \mbox{likelihood} & L(\{w_c\}) = \prod_{n=1}^{N} \prod_{c=1}^{C} \left(\frac{e^{z_c^{(n)}}}{\sum_{c'} e^{z_{c'}^{(n)}}} \right)^{\mathbb{I}(y^{(n)}=c)} \\ \mbox{log-likelihood} & \ell(\{w_c\}) = \sum_{n=1}^{N} \sum_{c=1}^{C} \mathbb{I}(y^{(n)}=c) (z_c^{(n)} - \log \sum_{c'} e^{z_{c'}^{(n)}}) \end{array}$$

one-hot encoding for labels
$$y^{(n)} \to [\mathbb{I}(y^{(n)} = 1), \dots, \mathbb{I}(y^{(n)} = C)]$$

 $z^{(n)} = [z_1^{(n)}, z_2^{(n)}, \dots, z_C^{(n)}], \quad z_c^{(n)} = w_c^{\top} x^{(n)}$

using this encoding from now on

log-likelihood
$$\ell(\{w_c\}) = \sum_{n=1}^N \left({y^{(n)}}^ op z^{(n)} - \log \sum_{c'} e^{z^{(n)}_{c'}}
ight)$$

Optimization

given the training data $\mathcal{D} = \{(x^{(n)}, y^{(n)})\}_n$ find the best model parameters $\{w_c\}_c$ by minimizing the cost (maximizing the likelihood of \mathcal{D})

$$J(\{w_c\}) = -\sum_{n=1}^N ({y^{(n)}}^ op z^{(n)} - \log \sum_{c'} e^{z^{(n)}_{c'}})$$
 where $z_c = w_c^ op x$

need to use gradient descent (for now calculate the gradient)

$$abla J(w) = [rac{\partial}{\partial w_{1,1}}J, \dots rac{\partial}{\partial w_{1,D}}J, \dots, rac{\partial}{\partial w_{C,D}}J]^ op$$
length $C imes D$

Gradient

