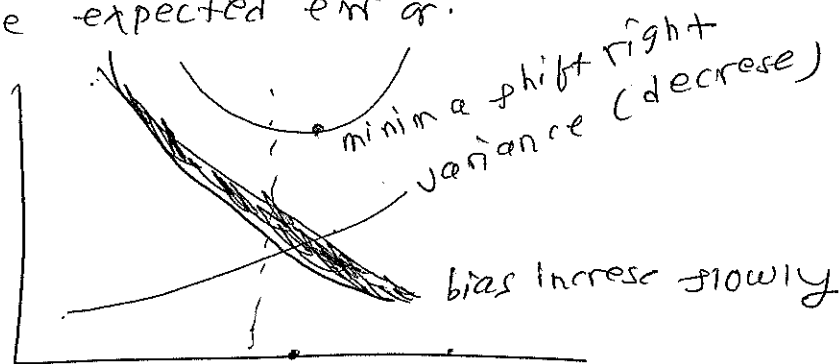


when training data increases, variance is reduced  
and slightly more complicated model minimizes  
the expected error.



① show that  $x^T A y$  is a valid kernel if  $A$  is sym. PSD.  
 $\Rightarrow A = Q^T \Lambda Q$  where  $\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$  and  $\lambda_1, \lambda_2, \dots, \lambda_n \geq 0$   
 define  $F = \text{diag}(\sqrt{\lambda_1}, \sqrt{\lambda_2}, \dots, \sqrt{\lambda_n})$  then  $\Lambda = F^T F$

$$\Rightarrow A = Q^T F^T F Q$$

$$\begin{aligned} \Rightarrow x^T A y &= x^T Q^T F^T F Q y \\ &= (F Q x)^T (F Q y) \\ &= \phi(x)^T \phi(y) \end{aligned}$$

$$\text{where } \phi(x) = F Q x$$

Aside:

$$(AB)^T = B^T A^T$$

$$(ABC)^T = C^T B^T A^T$$

$$\downarrow$$

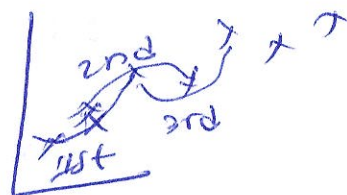
$$((AB)C)^T = C^T (AB)^T = C^T B^T A^T$$

## ⑥ cubic spline smoothing

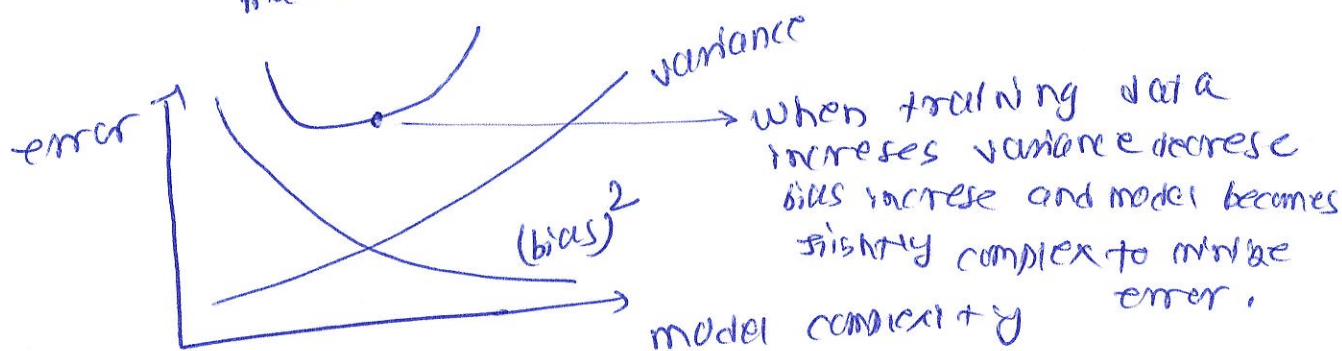
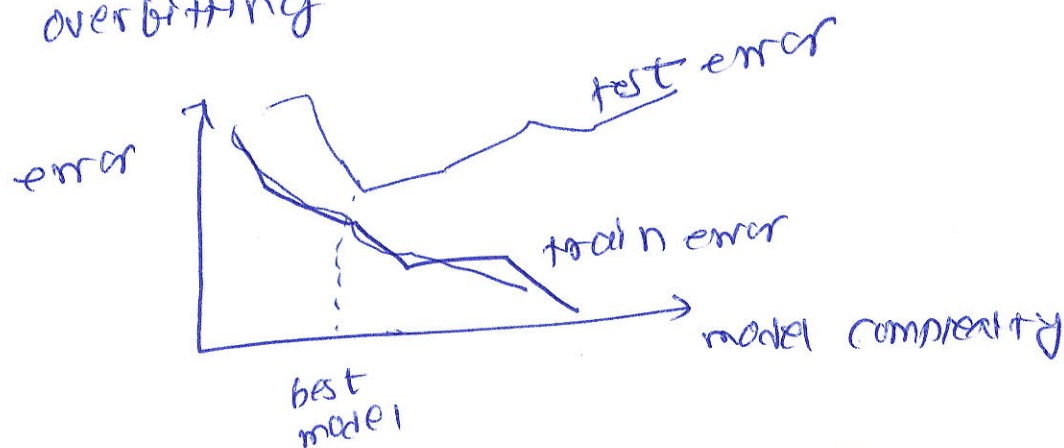
take 3 points  $x, x_{(i)}, x_{(i+1)}$

$$s_i(x) = a_i (x - x_{(i)})^3 + b_i (x - x_{(i)})^2 + c_i (x - x_{(i)}) + d_i$$

$$\forall x \in [x_{(i)}, x_{(i+1)}]$$



## ⑦ overfitting



# ① Likelihood function for logistic Regression

$$p(h_n | w) = h_n^n \cdot (1-h_n)^{1-n}$$

$$p(w) = \prod_{n=1}^N h_n^{t_n} (1-h_n)^{1-t_n}$$

$$\text{cost } E = \frac{1}{N} \cdot (-\ln p) = -\frac{1}{N} \sum_{n=1}^N [t_n \ln h_n + (1-t_n) \ln (1-h_n)]$$

$$\text{grad} = \nabla_w J = \frac{1}{N} \cdot \sum_{n=1}^N (h_n - t_n) \cdot x_n$$

initial  $z \leftarrow 0, \bar{w} = 0, \tau = 1$

for  $i = 1 : N$

$h_n = \text{sgn}(w^T x_n)$

if  $h_n \neq t_n$  then

$w = w + t_n x_n$

$\tau = \tau + 1$

$\bar{w} = \bar{w} + w$

$\tau = \tau + 1$

return  $\bar{w} / \tau$

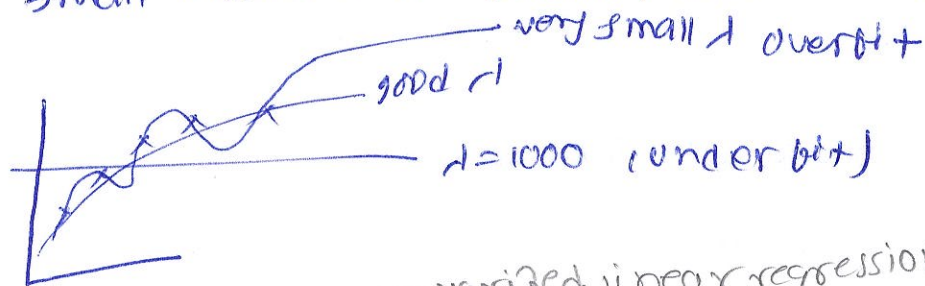
Test:  $\text{sgn}(w^T x)$

## ② perceptron criterion

$$\text{minimize } E_p(w) = - \sum_{n \in \text{misclassified}} t_n w^T x_n$$

①  $L_2 \rightarrow$  more penalize larger weights  
does not drive weights to zero  
weight vector is NOT sparse.

②  $\lambda \uparrow \Rightarrow$  high bias, low variance, underfit, simpler model  
drives weights closer to zero.  
small change in training data  $\Rightarrow$  big change in estimate



cost function for  $L_2$ -regularized linear regression

$$③ J = \frac{1}{2N} \sum_{n=1}^N (h_n - t_n)^2 + \frac{\lambda}{2} \|w\|^2$$

$N$ -examples  
in features

$$X = \begin{bmatrix} x_1^1 & x_1^2 & x_1^3 & x_1^m \\ x_2^1 & x_2^2 & x_2^3 & x_2^m \\ \vdots & \vdots & \vdots & \vdots \\ x_N^1 & x_N^2 & x_N^3 & x_N^m \end{bmatrix}$$

$$J(w) = \frac{1}{2N} \sum_{n=1}^N \left( w_0 + \sum_{j=1}^m w_j x_n^j - t_n \right)^2 + \frac{\lambda}{2} \sum_{j=1}^m w_j^2$$

(we omit  $w_0$  in regularization)

regularizing  $w_0$  = shifting origin of target

$\Rightarrow$  some change in all target values  $\Rightarrow$  similar change in estimates

④ Batch GD

stochastic GD

momentum GD  $\nabla J(w)$

$$v^{t+1} = \gamma \nabla J(w^t)$$

$$w^{t+1} = w^t - \eta \nabla J(w^t)$$

$$v^{t+1}$$

$\uparrow$

$\eta \approx 0.9$

for num-iter:

for num-iter:

for sample in data:

$$w^{t+1} = w^t - \eta \nabla J$$

$$w^{t+1} = w^t - \eta \nabla J$$

learning rate

Nesterov Accelerated Gradient

$$w^{t+1} = w^t - \eta \nabla J(w^t - \gamma \nabla J(w^t)) - \gamma \nabla J(w^t)$$

⑤ Gradient descent

$\eta$  = learning rate  
 $v$  = momentum  
 $\theta$  = hyperparameter  
 $w$  = weight vector at time stamp  $t$

$$v^{t+1} = \begin{cases} \eta \nabla J(w^t) & \text{vanilla} \\ \eta \nabla J(w^t) + v^t & \text{with momentum} \\ \eta \nabla J(w^t - v^t) + v^t & \text{Nesterov Accelerated Gradient (NAG)} \end{cases}$$

$$w^{t+1} = w^t - v^{t+1}$$

$\eta \approx 0.01$

⑥ Types of GD based on data used:

Batch GD  $\rightarrow$  use all data

Stochastic GD  $\rightarrow$  use only one example and update  $w$  after each iteration

Minibatch GD  $\rightarrow$  use constant number of examples and update  $w$  after each

$$J = \frac{1}{2N} \sum_{n=1}^N (h_n - t_n)^2$$

$$\nabla J = \frac{1}{N} \sum_{n=1}^N (h_n - t_n) \cdot x_n = \frac{1}{N} \cdot \text{np.sum}((h - t) \cdot x)$$

$$w = w - \eta \text{grad}$$

⑦ GD vs Normal eqn

$$w = w - \eta \text{grad}$$

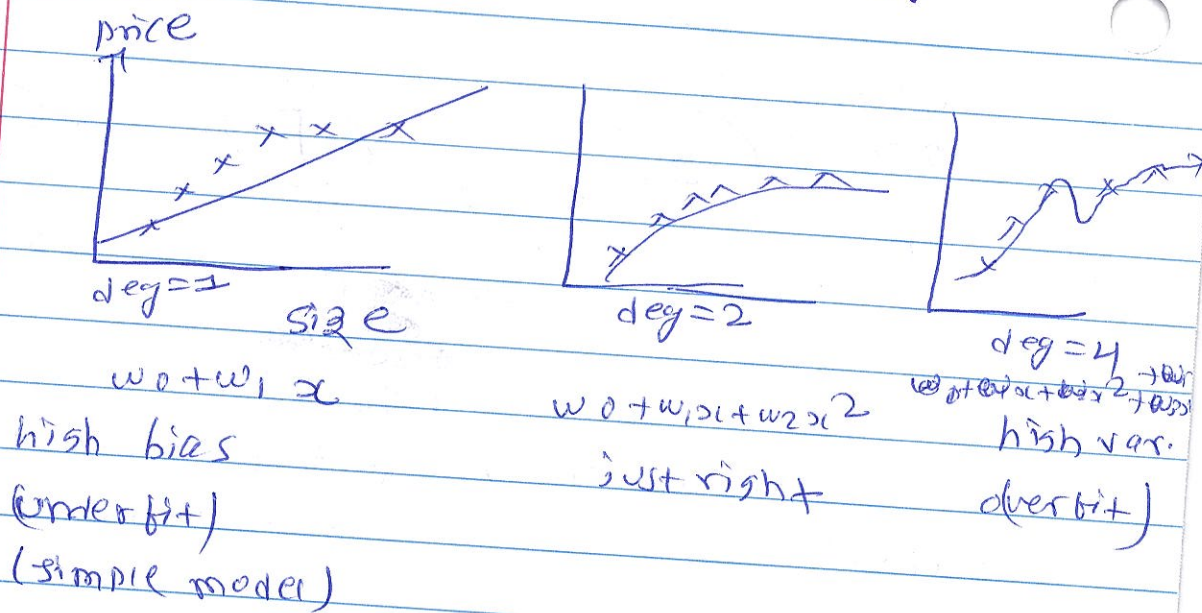
$$w = \underbrace{(X^T X)^{-1}}_{\text{pseudo inverse (pinv)}} X^T t$$

$$w = \underbrace{(X^T X + \lambda N I)^{-1}}_{L2} X^T t$$

for  $L2$  regularized  $I$  is identity matrix of shape  $X^T X$   
 $I[0][0] = 0$  for not to regularize bias term

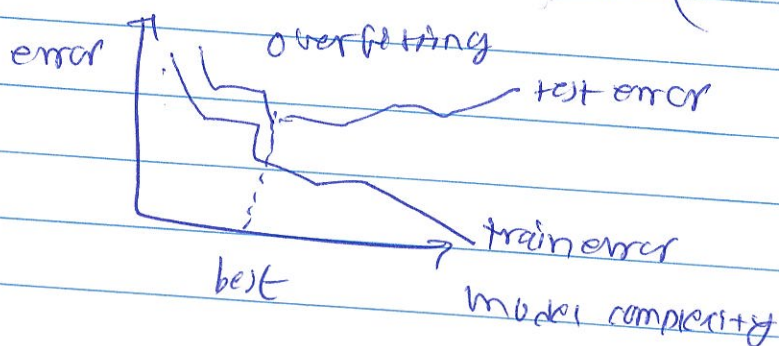
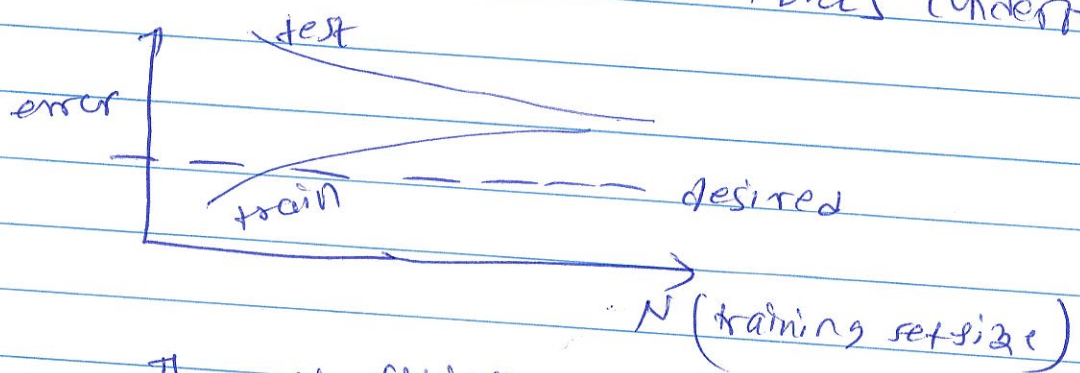


# Bias-variance Trade off

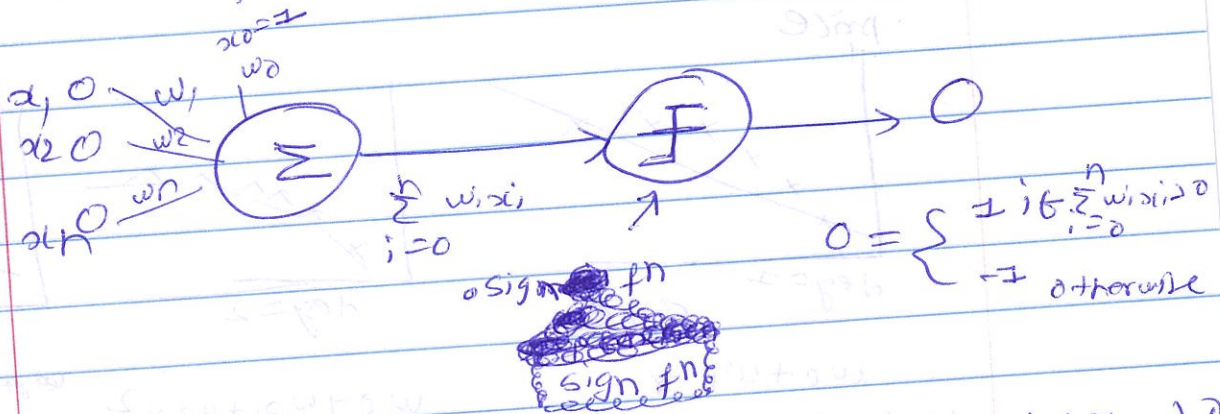


more training eg  
 smaller feature set  $\rightarrow$  bias high  $\sigma^2$  (underfitting)  
 larger feature set  $\rightarrow$  bias high  $\sigma^2$  (")  
 $\rightarrow$  bias high (underfitting)

## Learning curve for high bias (underfit)



# perceptron



$$\text{output} = \begin{cases} 1 & \text{if } \sum_{i=0}^n w_i x_i \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

$$= \begin{cases} 1 & \text{if } \vec{w} \cdot \vec{x} \geq 0 \\ -1 & \text{otherwise} \end{cases}$$



# ① cost for multivariate linear regression

$$J = \frac{1}{2N} \sum_{n=1}^N (h_n - t_n)^2$$

$$J = \frac{1}{2N} \sum_{n=1}^N \left( \sum_{j=0}^M w_j x_n^j - t_n \right)^2$$

$$h = w^T x \quad h = x^T @ w$$

$$h_n = \sum_{j=0}^M w_j x_n^j$$

$$h_1 = w_0 + w_1 x_1^{(1)} + w_2 x_1^{(2)} + w_M x_1^{(M)}$$

first row of sample

$$w = \begin{bmatrix} w_0 & w_1 & w_2 & w_3 \end{bmatrix} 1 \times 4$$

↑  
bias

$$h = (50, 4) \cdot (4, 1) = (80, 1)$$

$$X = \begin{bmatrix} 1 & 10 & 100 \\ 1 & 20 & 200 \\ 1 & 30 & 300 \\ 1 & 50 & 500 \end{bmatrix} 50 \times 4$$

$$t = \begin{bmatrix} 11K \\ 21K \\ 31K \\ 50K \end{bmatrix} 50 \times 1$$

2 features ~~age, birth year~~ <sup>age, birth size</sup>

$$e = h - t = \begin{bmatrix} h_1 - t_1 \\ h_2 - t_2 \\ \vdots \\ h_N - t_N \end{bmatrix} 50 \times 1$$

$$SSE = (h - t) \times x 2$$

$$MSE = \frac{(h - t) \times x 2}{N}$$

$$mse = np.mean(SSE)$$

$$rmse = np.sqrt(mse)$$

$$J = np.sum((h - t) \times x 2) / 2 / N$$

$$J = 0.5 * np.mean((h - t) \times x 2)$$

$$J = \frac{1}{2} \times mse \quad (\text{single float})$$

$$J = \frac{0.5}{len(t)} (h - t)^T (h - t)$$

$$(h - t)^T (h - t) = \begin{bmatrix} h_1 - t_1 & h_2 - t_2 & \dots \end{bmatrix} \begin{bmatrix} h_1 - t_1 \\ h_2 - t_2 \\ \vdots \\ h_N - t_N \end{bmatrix}$$

$$= \left[ \sum_n (h_n - t_n)^2 \right] 1 \times 1$$

## ② gradient of J

$$\nabla_w J = \nabla_w \left[ \frac{1}{2N} \sum_n (h - t)^2 \right] = \begin{bmatrix} \frac{\partial J}{\partial w_0} & \frac{\partial J}{\partial w_1} & \frac{\partial J}{\partial w_2} & \frac{\partial J}{\partial w_M} \end{bmatrix} 4 \times 1$$

$$= (h - t) \cdot T @ X^T$$

$$= (1, 50) \cdot (50, 4)$$

$$grad = \nabla_w J = (1, 4)$$

$$grad\_ols = (h - t) \cdot T @ X^T$$

$$\frac{\partial J}{\partial w_j} = \frac{1}{2N} \sum_{n=1}^N (x_n w_j - t_n) \cdot x_n$$

summation variables

$$\nabla_w J = \frac{1}{N} (h - t) \cdot X$$

we use matrix

①  $X = \text{design matrix}$   $\begin{bmatrix} & \end{bmatrix}_{N \times M}$   $N = \text{samples}$   
 $X_1 = \text{biased design matrix}$   $\begin{bmatrix} 1 & \end{bmatrix}_{N, M+1}$   $m = \text{features}$

$t = \begin{bmatrix} \end{bmatrix}_{N \times 1}$  column vector

$w = [w_0 \ w_1 \ w_2 \ w_m]_{1, M+1}$  row vector (2d array)

$w = \text{np.array}(w). \text{reshape}(1, \text{xt.shape}[1]) = (1, M+1)$

$h = X_1 @ w.T = (N, M+1) (M+1, 1) = (N, 1)$  same as  $t$

$e = h - t = N, 1 = 50, 1$

linear regression  $SSE = \sum_{n=1}^N (h-t)^2$

$MSE = \frac{1}{N} \sum_{n=1}^N (h-t)^2$

$J = \frac{1}{2N} \sum_{n=1}^N (h-t)^2 = \text{np.sum}((h-t) \times \times 2) / 2 / \text{float}(N)$

$RMSE = E_{\text{rms}} = \sqrt{MSE}$

Normal eqn:  $w = (X^T X)^{-1} X^T e$

moore penrose pseudo inverse of  $X$

$\text{np.linalg.pinv}(X)$

# multivariate linear

	bias	x		t
		geography	bed	2K
1		100	2	
1		200	3	3K
				5K
1		500	5	

t //

m features and one bias

$$X = \begin{bmatrix} x_0^{(1)} & x_1^{(1)} & \dots & x_m^{(1)} \\ x_0^{(2)} & x_1^{(2)} & \dots & x_m^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ x_0^{(N)} & x_1^{(N)} & \dots & x_m^{(N)} \end{bmatrix}_{N \times m+1} \quad \begin{bmatrix} t_1 \\ t_2 \\ t_3 \\ \vdots \\ t_n \end{bmatrix}_{N \times 1}$$

→ there are m features and one bias total (m+1)  
there are N samples for each feature

$$J = \frac{1}{2N} \|XW - t\|^2 = \frac{1}{2N} (XW - t)^T (XW - t)$$

$$\nabla_W J = \frac{1}{2} \frac{d}{d\alpha} (XW - t) \cdot 2(XW - t) \quad \frac{d}{d\alpha} (\alpha^T \alpha) = 2\alpha$$

α is a matrix

$$0 = X^T (XW - t)$$

$$X^T X W = X^T t$$

weights

$$W = (X^T X)^{-1} (X^T t)$$

$$\frac{d}{d\alpha} (A\alpha + b)^T (A\alpha + b) = 2A^T (A\alpha + b)$$

derivatives of ~~matrix~~ matrix products

# L1 vs L2 norm L2 more used

L2  $\rightarrow$  more penalty on large weights, but doesn't drive small weights to zero.

L1  $\rightarrow$  less penalty for large wt, but leads many weights ~~towards~~ to (or very very close) to zero. leading to weight vector to be sparse.

## Bias-Variance Tradeoff

(a) Ridge: 
$$\sum_{i=1}^N \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p \beta_j^2 = RSS + \lambda \sum_{j=1}^p \beta_j^2$$

shrinkage parameter  $\lambda$  which has to be tuned via validation set  $\rightarrow$  regularization strength

Ridge  $\left\{ \begin{array}{l} \lambda \uparrow \Rightarrow \text{var} \downarrow \text{bias} \uparrow \\ \text{meaning: small change in training data} \\ \text{big change in parameter estimates} \\ \text{effect will increase with no. of parameters} \end{array} \right.$

(b) LASSO  $\beta_j^2$  vs  $|\beta_j|$

Ridge shrinks wts but do not make 0  
but LASSO makes them zero, makes less no. of wts.



# Logistic Regression → classification (binary)

linear regression  $h = w^T x$

logistic regression  $h = \sigma(w^T x) = \frac{1}{1 + e^{-w^T x}} = h(x) = y(x)$

likelihood function  $p(t|w) = \prod_n h_n^{t_n} (1 - h_n)^{1 - t_n}$

→ log likelihood,  $E$  or  $J = -\ln p(t|w)$

$$= -\sum_{n=1}^N [t_n \ln h_n + (1 - t_n) \ln (1 - h_n)]$$

∴ cost or Error or loss function

$$E_D = -\sum_n [t_n \ln h_n + (1 - t_n) \ln (1 - h_n)]$$

where  $t_n \in \{0, 1\}$  NOT  $\{-1, 1\}$

add regularization,

$$E_w = \frac{\lambda}{2} w^T w$$

→ then  $L_2$ -regularized logistic regression cost function

~~cost~~  $E = -\frac{1}{N} \sum_n [t_n \ln h_n + (1 - t_n) \ln (1 - h_n)] + \frac{\lambda}{2} w^T w$

$\lambda$  is regularized

$$h(x) = \sigma(w^T x) = \frac{1}{1 + e^{-w^T x}}$$

$$h_n = h(x_n) = \sigma(x_n) = \frac{1}{1 + e^{-w^T x_n}} \text{ is sigmoid fn}$$

# Softmax Regression (classification / multiclass)

training set:  $(x_1, t_1), (x_2, t_2), \dots, (x_n, t_n)$

$$X = [x_1, x_2, \dots, x_n]$$

$$t_1, t_2, \dots, t_n \in \{1, 2, \dots, K\}$$

$$w_k = [w_{k0}, w_{k1}, w_{k2}, \dots]^T$$

one weight vector per class,

$$p(c_k | x) = \frac{e^{w_k^T x}}{\sum_j e^{w_j^T x}}$$

prob that  
nth example  $x_n$   
has class  $t_n$

$$p(t_n | x_n) = \frac{e^{w_{t_n}^T x_n}}{\sum_j e^{w_j^T x_n}}$$

likelihood is the joint probability of all classes

$$L(w) = \prod_{n=1}^N p(t_n | x_n)$$

cost is the -ve log likelihood,

$$E_D(w) = -\frac{1}{N} \ln L(w)$$

$$= -\frac{1}{N} \ln \prod_{n=1}^N p(t_n | x_n)$$

$$= -\frac{1}{N} \sum_n \ln p(t_n | x_n)$$

$$= -\frac{1}{N} \sum_n \ln \frac{e^{w_{t_n}^T x_n}}{\sum_j e^{w_j^T x_n}}$$

$w_{t_n}$  is weight vector for  
nth example

$w_k$  is weight vector for  
all the examples belonging  
to class  $k$ .

$$E_D(w) = -\frac{1}{N} \sum_{n=1}^N \sum_{k=1}^K \delta_k(t_n) \ln \left( \frac{e^{w_k^T x_n}}{\sum_j e^{w_j^T x_n}} \right)$$

## MLE

sample:  $x_1, x_2, \dots, x_N$

prob of  $x$  is  
 $b(x, w)$

joint prob of  $x_1, x_2, \dots, x_N$  is called

likelihood  $L(w)$

$$L(w) = p(x_1, x_2, \dots, x_N | w) = b(x_1, w) \dots b(x_N, w)$$

$$L(w) = \prod_{n=1}^N b(x_n, w)$$

# MAP solution for LR

$$p(t|w) = \prod_{n=1}^N h_n^{t_n} (1-h_n)^{1-t_n}$$

$$p(w) = \left(\frac{1}{2\pi}\right)^{\frac{m+1}{2}} e^{-\frac{1}{2} w^T w} \propto e^{-\frac{1}{2} w^T w}$$

$$p(w|t) = \frac{p(t|w) p(w)}{p(t)} \propto p(t|w) p(w)$$

$$\text{now, } \hat{w}_{\text{MAP}} = \underset{w}{\operatorname{argmax}} p(w|t)$$

$$= \underset{w}{\operatorname{argmax}} p(t|w) p(w)$$

$$= \underset{w}{\operatorname{argmin}} -\ln p(t|w) -\ln p(w)$$

$$= \underset{w}{\operatorname{argmin}} -\ln \prod_n h_n^{t_n} (1-h_n)^{1-t_n} -\ln e^{-\frac{1}{2} w^T w}$$

$$= \underset{w}{\operatorname{argmin}} -\sum_n \ln(h_n^{t_n} (1-h_n)^{1-t_n}) + \frac{1}{2} w^T w$$

$$= \underset{w}{\operatorname{argmin}} -\sum_n \underbrace{[t_n \ln h_n + (1-t_n) \ln (1-h_n)]}_{E_D(w)} + \underbrace{\frac{1}{2} w^T w}_{E_w(w)}$$

$$\boxed{\hat{w}_{\text{MAP}} = \underset{w}{\operatorname{argmin}} E_D(w) + E_w(w)}$$

ans

# Softmax Regression

$$P(k|x) = \frac{e^{w_k^T x}}{\sum_j e^{w_j^T x}}$$

$$P(t_n|x_n) = \frac{e^{w_{t_n}^T x_n}}{\sum_{j=1}^K e^{w_j^T x_n}}$$

note

$$h(w) = \frac{1}{\sum_{k=1}^K e^{w_k^T x}} \begin{bmatrix} e^{w_1^T x} \\ e^{w_2^T x} \\ \vdots \\ e^{w_K^T x} \end{bmatrix}$$

$$E_D(w) = \frac{1}{N} \sum_n -\ln \prod_n P(t_n|x_n)$$

$$= -\frac{1}{N} \sum_n \ln \frac{e^{w_{t_n}^T x_n}}{\sum_j e^{w_j^T x_n}}$$

$$= -\frac{1}{N} \sum_n \ln \frac{e^{w_{t_n}^T x_n}}{\sum_j e^{w_j^T x_n}}$$

$$E_D(w) = -\frac{1}{N} \sum_n \sum_k \delta_k(t_n) \ln \left( \frac{e^{w_k^T x_n}}{\sum_k e^{w_k^T x_n}} \right)$$

Regularizer

$$+ \frac{\alpha}{2} \sum_{k=1}^K w_k^T w_k$$

$$\frac{\partial E_D}{\partial w_k} = \alpha w_k$$

(no summation)

Let,  $E_n = \sum_k \delta_k(t_n) w_k^T x_n - \sum_k \delta_k(t_n) \ln \left( \sum_k e^{w_k^T x_n} \right)$

$$\frac{\partial E_n}{\partial w_j} = \delta_j(t_n) x_n - \delta_j(t_n) \cdot \frac{1}{\sum_j e^{w_j^T x_n}} \cdot e^{w_j^T x_n} \cdot \delta_j(t_n) x_n$$

$$\frac{\partial E_n}{\partial w_k} = \delta_k(t_n) x_n - \frac{e^{w_k^T x_n}}{\sum_k e^{w_k^T x_n}} \cdot x_n$$

$$\frac{\partial E_n}{\partial w_k} = [\delta_k(t_n) - p(k|x_n)] x_n$$

$$\Rightarrow \frac{\partial E_D(w)}{\partial w_k} = -\frac{1}{N} \sum_n [\delta_k(t_n) - p(k|x_n)] x_n$$

there are K such equations for each classes.

prevent overflow

$$p(k|x) = \frac{e^{w_k^T x}}{\sum_k e^{w_k^T x}}$$

$$c = \max_{1 \leq k \leq K} w_k^T x$$

$$p(k|x) = \frac{e^{(w_k^T x - c)}}{\sum_k \exp(w_k^T x - c)}$$



## ① MLE VS MAP

MLE = maximize  $P(\text{data} | \text{params})$  by searching over parameters

MAP = maximize  $P(\text{param} | \text{data})$  by searching over params and accounting for prior over params.

MLE  $\rightarrow$  finds  $w$  by maximizing likelihood  $P(w | \text{data})$

MAP  $\rightarrow$  maximizes the posterior prob  $P(w | \text{data})$

② Logistic Regression maps inputs to discrete outputs  
" " continuous ".

## ③ PCA & feature selection

Similarity : reduce the dimension of data  
↓  
Reference : feature selection finds a subset of features  
PCA produces a smaller <sup>new</sup> set

## perceptron

perceptron criterion:  $w^T x_n \geq 0$  for  $t_n = +1$   
 $w^T x_n < 0$  for  $t_n = -1$

want:  $t_n w^T x_n \geq 0$  for all patterns  
minimize:  $-w^T x_n t_n$  for all misclassified patterns  $n$

$$\boxed{E_P(w) = - \sum_{n \in M} w^T x_n t_n}$$

perceptron ↑ mistakes (misclassified)

## Binary perceptron

initialize  $\vec{w} = 0$

for  $n = 1, \dots, N$

$h_n = \text{sgn}(w^T x_n)$

if  $h_n \neq t_n$  then

$w = w + t_n x_n$

} repeat until convergence  
or  
given number of epochs

① why kernels are symmetric?

Inner products are symmetric by definitions, so, therefore if the kernel function represent an inner product in some Hilbert space, then the kernel function must be symmetric as well.

$$\begin{aligned}K(x, y) &= \langle \phi(x), \phi(y) \rangle \\&= \langle \phi(y), \phi(x) \rangle \quad (\because \text{property of inner product}) \\&= K(y, x)\end{aligned}$$

② show  $K(x, z) = \alpha^T A^T A z$  is a valid kernel.

Let  $\phi(x) = Ax$ ,

then,

$$\begin{aligned}\langle \phi(x), \phi(z) \rangle &= \phi(x)^T \phi(z) \\&= (Ax)^T (Az) \\&= \alpha^T A^T A z \\&= K(x, z)\end{aligned}$$

$\therefore K(x, z)$  is an inner product in some Hilbert space.

③  $AA^T$  is PSD matrix (indirect)

proof let,  $A \in \mathbb{R}^{m \times n}$

$\lambda$  = eigenvalue of  $AA^T$

$q$  = eigenvector of  $\lambda$

$$(AA^T)q = \lambda q \quad (\text{premultiply by } q^T)$$

$$\Rightarrow q^T AA^T q = q^T \lambda q$$

$$\therefore AA^T = A^T A$$

$$\begin{aligned} \Rightarrow \lambda &= \frac{q^T AA^T q}{q^T q} = \frac{q^T A^T A q}{q^T q} \\ &= \frac{z^T z}{q^T q} \quad \text{where } z = A^T q \end{aligned}$$

$$\text{here } z^T z \geq 0$$

$$q^T q \geq 0$$

$\therefore \lambda \geq 0$  so  $AA^T$  is PSD.



## MLE vs MAP

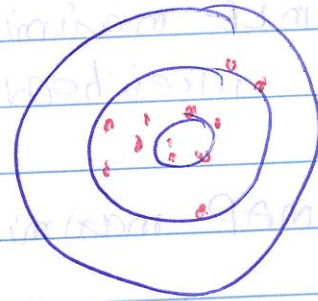
- MLE maximizes the  $\ln p(D|w)$   
likelihood of model  $w$  w.r.t. data  $D$
  - MAP maximizes the  $\ln p(w|D)$   
likelihood of data  $D$  w.r.t. model  $w$
- moreover, MAP additionally uses priors.

QAM 2v - FJM

Low variance

High variance

Low Bias



High Bias

