

# Machine Learning

## The Motivations & Applications of Machine Learning -

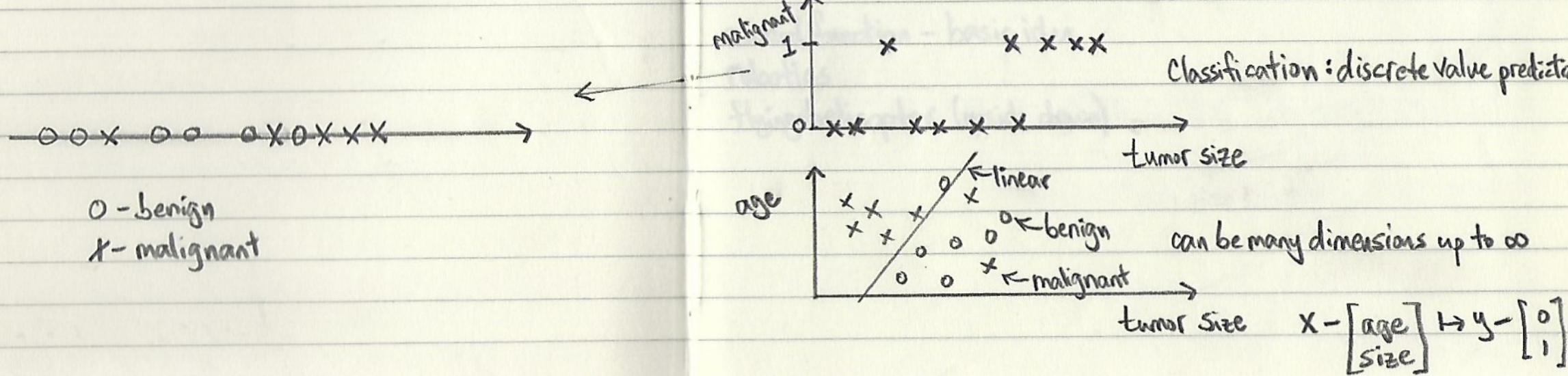
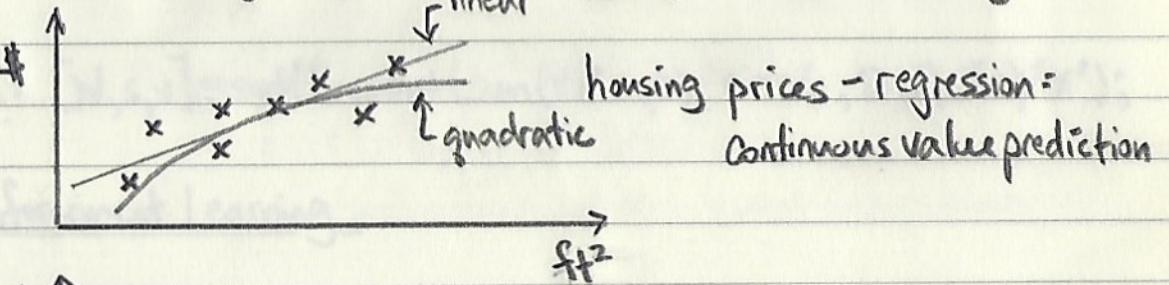
optical character recognition - US post zipcode  
 flying helicopter  
 mining databases - medical records  
 identifying fraudulent transactions  
 personalized recommendations  
 understanding human genome

Machine Learning definition:

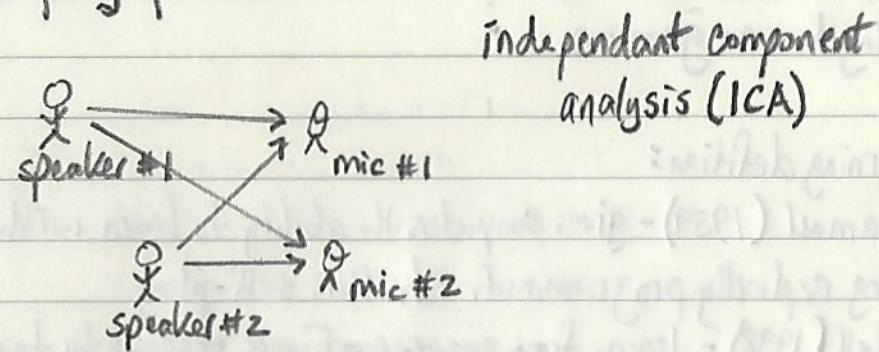
Arthur Samuel (1959) - gives computers the ability to learn without being explicitly programmed. Checkers self-play

Tom Mitchell (1998) - learns from experience  $E$  with respect to task  $T$  and some performance measure  $P$ , if its performance on  $T$ , as measured by  $P$ , improves with experience  $E$ .

#1 Supervised Learning: learn  $x \rightarrow y$  mappings    #2 Deep Learning.



Cocktail party problem:



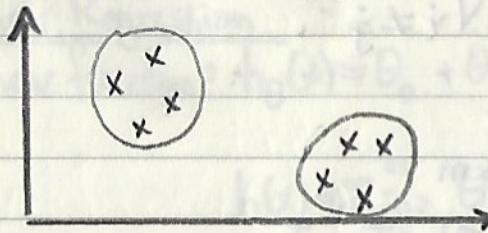
### #3 Learning Theory - Practical ML advice

When can you guarantee that a learning algorithm will work?

What algorithms will approximate which functions?

How much data is needed to train?

### #4 Unsupervised Learning - only input X



clustering -  
discover structure in the data

- gene analysis
- image clustering - grouping pixels together
- social network analysis
- market segmentation
- astronomical data analysis
- organize computing clusters
- grouping news articles

$$\text{ICA: } [W, s, v] = \text{svd}(\text{repmat}(\text{sum}(X \cdot X', 1), \text{size}(X, 1), 1) \cdot X' \cdot X');$$

### #5 Reinforcement Learning -

reward function - basic idea

robotics

flying helicopter (upside down)

Linear Algebra

vector  $v \in \mathbb{R}^n$  - n-dim real space

matrix  $A \in \mathbb{R}^{m \times n}$   $\begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{m1} & \dots & a_{mn} \end{bmatrix}$  m row vectors  
n col vectors

identity  $I^n \begin{bmatrix} 1 & 0 \\ \vdots & \ddots \\ 0 & 1 \end{bmatrix}$

diagonal  $D^n \begin{bmatrix} * & 0 \\ \vdots & \ddots \\ 0 & * \end{bmatrix} D_{ij} = 0 \forall i \neq j$

transpose  $A \in \mathbb{R}^{m \times n}$ ,  $A^T \in \mathbb{R}^{n \times m}$   
 $A_{ij} = (A^T)_{ji}$   $(A+B)^T = A^T + B^T$

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

inner product  $\hat{v} \in \mathbb{R}^n$   $u \in \mathbb{R}^n$

$$v^T u = \langle v, u \rangle = \sum_{i=1}^n v_i u_i \in \mathbb{R}$$

$$\mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$$

outer product  $\hat{v} \in \mathbb{R}^m$   $u \in \mathbb{R}^n$

$$v \cdot u^T \in \mathbb{R}^{m \times n}$$

$$(v \cdot u^T)_{ij} = v_i u_j$$

CS229:S1

CS229:N1

Supervised Learning

$(x^{(i)}, y^{(i)})$  - training example  $\{(x^{(i)}, y^{(i)}) : i=1\dots m\}$  - training set

$\uparrow$  input  $\uparrow$  target

$$X = Y = \mathbb{R}$$

learn function  $h: X \mapsto Y$  so  $h(x)$  is a good predictor  
 $\uparrow$  hypothesis

Linear Regression

Linear function:  $h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$ ,  $\theta_i$  - parameters, weights  
 $\uparrow$  features

$$h(x) = \sum_{i=0}^n \theta_i x_i = \theta^T x, x_0 = 1 - \text{intercept term}$$

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 - \text{cost function}$$

$\uparrow$  ordinary least squares

choose  $\theta$  to minimize  $J(\theta)$

gradient descent:  $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta), j=0\dots n$

LMS update:  $\theta_j := \theta_j - \alpha (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)}$  - single training example  
repeat until convergence {

$$\theta_j := \theta_j - \alpha \sum_{i=1}^m (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)}, \text{for every } j \}$$

$\uparrow$  learning rate  
 $\uparrow$  batch gradient descent

$J$  is a convex quadratic function - gradient descent always converges

loop {

for  $i=1$  to  $m$  {

$$\theta_j := \theta_j + \alpha (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)}, \text{for every } j \}$$

$\uparrow$  stochastic gradient descent

$\uparrow$  may never converge

CS229:51

CS229=N1

matrix-vector:  $A \in \mathbb{R}^{m \times n}, x \in \mathbb{R}^n$ 

$$1) A = [a_1, a_2, \dots, a_n] \quad a_i \in \mathbb{R}^m$$

$$y = Ax = a_1x_1 + a_2x_2 + \dots + a_nx_n$$

↳ linear combination of columns of A

$$2) A = [a_1^T, \dots, a_m^T]^T \quad a_i \in \mathbb{R}^n$$

$$y = Ax = \begin{bmatrix} a_1^T \\ \vdots \\ a_m^T \end{bmatrix} x \in \mathbb{R}^{m \times 1}$$

matrix-matrix:  $A \in \mathbb{R}^{m \times k}, B \in \mathbb{R}^{k \times n}$ 

$$A \cdot B \in \mathbb{R}^{m \times n}$$

$$A = \begin{bmatrix} & & \\ & \ddots & \\ & & \end{bmatrix} \downarrow m \quad B = \begin{bmatrix} & & \\ & \ddots & \\ & & \end{bmatrix} \downarrow n$$

$\overbrace{\quad}^k \quad \overbrace{\quad}^n$

operations &amp; properties:

$$\text{trace} - \text{tr}(A) = \sum_{i=1}^n A_{ii}, \mathbb{R}^{n \times n} \rightarrow \mathbb{R} \text{ only square matrix}$$

range & null space:  $A \in \mathbb{R}^{m \times n}$ 

$$\text{range of } A = R(A) = \{v \in \mathbb{R}^m : v = A \cdot x, x \in \mathbb{R}^n\}$$

projection: projection of vector  $y \in \mathbb{R}^m$  onto span of  $\{x_1, \dots, x_n\}$  is

$$v = \text{Proj}_X(y; X) \quad \text{vector } v \text{ s.t. } v \text{ is as close as possible to } y$$

$$= \arg \min_{v \in R(X)} \|v - y\|_2$$

$$= X(X^T X)^{-1} X^T y$$

$$\text{null space} = \{x \in \mathbb{R}^n, Ax = 0\}$$

The Normal Equationsfor function  $f: \mathbb{R}^{n \times d} \mapsto \mathbb{R}$ ,  $\nabla_A f(A)$  ↴  $\uparrow$   $m \times n$  matrix

$$\left[ \begin{array}{c} \frac{\partial f}{\partial A_{11}} \cdots \frac{\partial f}{\partial A_{1d}} \\ \vdots \\ \frac{\partial f}{\partial A_{n1}} \cdots \frac{\partial f}{\partial A_{nd}} \end{array} \right] \downarrow n$$

$$\text{trace operator: } \text{tr } A = \sum_{i=1}^d A_{ii} \quad \longleftrightarrow d \quad \text{sum of diagonal entries of } n \text{-by-} n A$$

$$\text{tr } AB = \text{tr } BA$$

$$\text{tr } A = \text{tr } A^T$$

$$\text{tr } (A+B) = \text{tr } A + \text{tr } B \quad \text{cyclic permutation property}$$

$$\text{tr } aA = a \text{tr } A \quad B \text{ fixed}$$

$$\nabla_A \text{tr } AB = B^T$$

$$\nabla_A^T f(A) = (\nabla_A f(A))^T$$

$$\nabla_A \text{tr } AB A^T C = CAB + C^T A B^T$$

$$\nabla_A |\text{A}| = |\text{A}|(\text{A}^{-1})^T \quad \text{non-singular square A, } |\text{A}| \text{ is determinant}$$

$$X = \begin{bmatrix} -(x^{(1)})^T \\ \vdots \\ -(x^{(d)})^T \end{bmatrix} \downarrow n, \vec{y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(d)} \end{bmatrix}, X\theta - \vec{y} = \begin{bmatrix} h_\theta(x^{(1)}) - y^{(1)} \\ \vdots \\ h_\theta(x^{(d)}) - y^{(d)} \end{bmatrix}$$

$\longleftrightarrow d \longrightarrow$

$$\frac{1}{2} (X\theta - \vec{y})^T (X\theta - \vec{y}) = \frac{1}{2} \sum_{i=1}^n (h_\theta(x^{(i)}) - y^{(i)})^2 = J(\theta)$$

recall  $\vec{z}^T \vec{z} = \sum_i z_i^2$

$$\begin{aligned} \nabla_\theta J(\theta) &= \nabla_\theta \frac{1}{2} (X\theta - \vec{y})^T (X\theta - \vec{y}) \\ &= X^T X \theta - X^T \vec{y} \\ &= 0 \end{aligned}$$

$$\therefore X^T X \theta = X^T \vec{y}$$

$$\theta = (X^T X)^{-1} X^T \vec{y} \quad \text{- closed form}$$

determinant:  $A \in \mathbb{R}^{n \times n}$   
 $\det(A) \rightarrow \mathbb{R}$

$$|A| = |A^T|$$

$$|A \cdot B| = |A| \cdot |B|$$

$|A| = 0$ ,  $A$  is not full rank

$|A^{-1}| = \frac{1}{|A|}$  if  $A$  is invertible

$\text{rank}(A)$ : size of largest subset of columns of  $A$  that are linearly independent

$$A \in \mathbb{R}^{n \times n}$$

eigenvalue of  $A$  is  $\lambda \in \mathbb{C}$  (complex number)

$u$  is eigenvector of  $A$  if  $Au = \lambda u$  ( $u \neq 0$ )

$$Au - \lambda u = 0 \Leftrightarrow (A - \lambda I)u = 0$$

↑

$$\det(A - \lambda I) = 0$$

1)  $\forall A \in \mathbb{R}^{n \times n}$ ,  $n$  eigenvalues (not necessarily unique)

$$\det(I - \lambda I) = (\lambda - 1)^n$$

2)  $\forall u$  eigenvector of  $A$

$$\exists \lambda \text{ s.t. } Au = \lambda u$$

$\forall c \neq 0$ ,  $cu$  also an eigenvector  $A(cu) = \lambda(cu)$

$\lambda$  eigenvalue  $u_1, u_2$  eigenvectors,  $Au_1 = \lambda u_1$ ,  $Au_2 = \lambda u_2$

$\forall c_1, c_2 - c_1 u_1 + c_2 u_2$  is still eigenvector of  $\lambda$

$$3) \text{tr}(A) = \sum_{i=1}^n \lambda_i$$

$$\det(A) = \prod_{i=1}^n \lambda_i$$

$$\text{rank}(A) = \text{nnz}(\{\lambda_i\}_{i=1}^n)$$

### Probabilistic Interpretation

$$y^{(i)} = \theta^T x^{(i)} + \varepsilon^{(i)}, \varepsilon^{(i)} \text{ is error term, IID (independent, identically distributed)}$$

$$\varepsilon^{(i)} \sim N(0, \sigma^2), p(\varepsilon^{(i)}) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(\varepsilon^{(i)})^2}{2\sigma^2}\right)$$

$\uparrow$  density of  $\varepsilon^{(i)}$

$$p(y^{(i)} | x^{(i)}; \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)$$

$\downarrow$  parameterized by  $\theta$ , not a random variable

$$\text{distribution of } y^{(i)} = y^{(i)} | x^{(i)}; \theta \sim N(\theta^T x^{(i)}, \sigma^2)$$

$$\text{likelihood function } L(\theta) = L(\theta; X, \bar{y}) = p(\bar{y} | X; \theta)$$

$$= \prod_{i=1}^m p(y^{(i)} | x^{(i)}; \theta)$$

$$= \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)$$

$$\log \text{likelihood } l(\theta) = \log L(\theta) = m \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{2\sigma^2} \cdot \frac{1}{2} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2$$

maximizing  $l(\theta)$  means minimizing  $\frac{1}{2} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2$   
 independent of  $\sigma$   $\uparrow$  same as  $J(\theta)$  least-squares

### Locally Weighted Linear Regression (LWR)

1. fit  $\theta$  to minimize  $\sum_i w^{(i)} (y^{(i)} - \theta^T x^{(i)})^2$

2. output  $\theta^T x$   $\uparrow$  non-negative weights

standard choice for  $w^{(i)} = \exp\left(-\frac{(x^{(i)} - x)^2}{2\tau^2}\right)$ ,  $\tau$  is bandwidth parameter

non-parametric algorithm - need entire training set to make predictions

## Classification and Logistic Regression

values of  $y$  take on small number of discrete values

$$h_{\theta}(x) = g(\theta^T x) = \frac{1}{1+e^{-\theta^T x}}, \quad g(z) = \frac{1}{1+e^{-z}}, \quad g'(z) = g(z)(1-g(z))$$

$\uparrow$  sigmoid

$$\begin{aligned} \text{assume } P(y=1|x; \theta) &= h_{\theta}(x) \\ P(y=0|x; \theta) &= 1 - h_{\theta}(x) \end{aligned} \quad \left. \begin{aligned} p(y|x; \theta) &= (h_{\theta}(x))^y (1-h_{\theta}(x))^{1-y} \\ &\uparrow \text{binary classification} \end{aligned} \right\}$$

$$\begin{aligned} \text{likelihood of parameters } L(\theta) &= p(\vec{y}|x; \theta) = \prod_{i=1}^m (h_{\theta}(x^{(i)}))^{y^{(i)}} (1-h_{\theta}(x^{(i)}))^{1-y^{(i)}} \\ l(\theta) &= \log L(\theta) = \sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1-y^{(i)}) \log (1-h_{\theta}(x^{(i)})) \end{aligned} \quad \left. \begin{aligned} &\uparrow \text{maximize} \end{aligned} \right.$$

$$\text{gradient ascent } \theta := \theta + \alpha \nabla_{\theta} l(\theta)$$

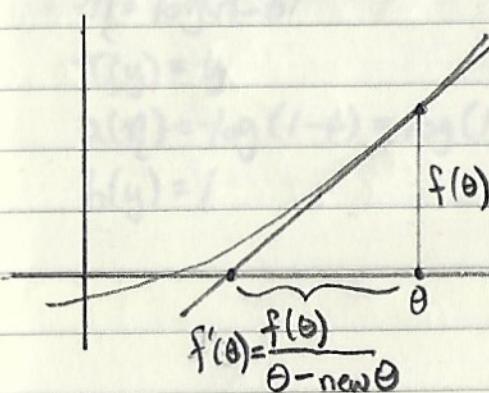
$\uparrow$  + because maximizing log likelihood

$$\frac{\partial}{\partial \theta_j} l(\theta) = (y - h_{\theta}(x)) x_j$$

$$\text{stochastic gradient ascent } \theta_j := \theta_j + \alpha (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$$

Newton's method: to find  $\theta$  such that  $f(\theta) = 0$ , use update rule

$$\theta := \theta - \frac{f(\theta)}{f'(\theta)}, \quad \theta \in \mathbb{R}$$



To maximize  $l(\theta)$ , find  $l'(\theta) = 0$   
using Newton's method:  
 $\theta := \theta - \frac{l''(\theta)}{l'(\theta)}$

$$\text{Gaussian distribution} - p(y; \mu) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(y-\mu)^2\right)$$

$$= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}y^2\right) \cdot \exp(\mu y - \frac{1}{2}\mu^2)$$

$\sigma^2 = 1$  since  $\sigma^2$  has no effect on  $\Theta$  or  $h_\theta(x)$

$$\eta = \mu, T(y) = y, a(\eta) = \eta^2/2 = \mu^2/2, b(y) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{y^2}{2}\right)$$

Constructing GLMs - predict random  $y$  as a function of  $x$   
 three assumptions for GLM:

1.  $y|x; \Theta \sim \text{Exponential Family}(\eta)$
2. given  $x$ , predict expected value of  $T(y)$  given  $x$   
 if  $T(y) = y$ ,  $h(x) = E[y|x]$
3.  $\eta = \Theta^T x$  -  $\eta$  and inputs are related linearly

ordinary least squares -  $h(x) = E[y|x; \Theta] = \mu = \eta = \Theta^T x$   
 $y|x; \Theta \sim N(\mu, \sigma^2)$  - Gaussian

logistic regression -  $h(x) = E[y|x; \Theta] = \phi = \frac{1}{1+e^{-\eta}} = \frac{1}{1+e^{-\Theta^T x}}$   
 $y|x; \Theta \sim \text{Bernoulli}(\phi)$

canonical response function -  $g(\eta) = E[T(y); \eta]$   
 canonical link function -  $g^{-1}$

softmax regression -  $y \in \{1, 2, \dots, K\}$ ,  $K$  classes

parameterize multinomial with  $\phi_1, \dots, \phi_{K-1}$

$$\phi_i = p(y=i; \phi) \text{ and } p(y=K; \phi) = 1 - \sum_{i=1}^{K-1} \phi_i$$

$\phi_K = 1 - \sum_{i=1}^{K-1} \phi_i$  - not a parameter since fully specified  
 by  $\phi_1, \dots, \phi_{K-1}$

1/21/18

$$(x_i - \mu)^T (x_j - \mu) = (x_i - \mu)^T q = \text{middle term in variance}$$

$$(x_i - \mu)^T (x_j - \mu) = (x_i - \mu)^T (x_j - \mu)$$

$$(x_i - \mu)^T (x_j - \mu) = (x_i - \mu)^T (x_j - \mu)$$

$$(x_i - \mu)^T (x_j - \mu) = (x_i - \mu)^T (x_j - \mu)$$

$$(x_i - \mu)^T (x_j - \mu) = (x_i - \mu)^T (x_j - \mu)$$

$$(x_i - \mu)^T (x_j - \mu) = (x_i - \mu)^T (x_j - \mu)$$

$$(x_i - \mu)^T (x_j - \mu) = (x_i - \mu)^T (x_j - \mu)$$

$$(x_i - \mu)^T (x_j - \mu) = (x_i - \mu)^T (x_j - \mu)$$

$$(x_i - \mu)^T (x_j - \mu) = (x_i - \mu)^T (x_j - \mu)$$

$$(x_i - \mu)^T (x_j - \mu) = (x_i - \mu)^T (x_j - \mu)$$

$$(x_i - \mu)^T (x_j - \mu) = (x_i - \mu)^T (x_j - \mu)$$

$$(x_i - \mu)^T (x_j - \mu) = (x_i - \mu)^T (x_j - \mu)$$

$$(x_i - \mu)^T (x_j - \mu) = (x_i - \mu)^T (x_j - \mu)$$

$$(x_i - \mu)^T (x_j - \mu) = (x_i - \mu)^T (x_j - \mu)$$

$$(x_i - \mu)^T (x_j - \mu) = (x_i - \mu)^T (x_j - \mu)$$

$$(x_i - \mu)^T (x_j - \mu) = (x_i - \mu)^T (x_j - \mu)$$

$$(x_i - \mu)^T (x_j - \mu) = (x_i - \mu)^T (x_j - \mu)$$

$$(x_i - \mu)^T (x_j - \mu) = (x_i - \mu)^T (x_j - \mu)$$

$$(x_i - \mu)^T (x_j - \mu) = (x_i - \mu)^T (x_j - \mu)$$

$$(x_i - \mu)^T (x_j - \mu) = (x_i - \mu)^T (x_j - \mu)$$

$$(x_i - \mu)^T (x_j - \mu) = (x_i - \mu)^T (x_j - \mu)$$

$$(x_i - \mu)^T (x_j - \mu) = (x_i - \mu)^T (x_j - \mu)$$

$$(x_i - \mu)^T (x_j - \mu) = (x_i - \mu)^T (x_j - \mu)$$

$$(x_i - \mu)^T (x_j - \mu) = (x_i - \mu)^T (x_j - \mu)$$

$$(x_i - \mu)^T (x_j - \mu) = (x_i - \mu)^T (x_j - \mu)$$

$$(x_i - \mu)^T (x_j - \mu) = (x_i - \mu)^T (x_j - \mu)$$

$$(x_i - \mu)^T (x_j - \mu) = (x_i - \mu)^T (x_j - \mu)$$

$$(x_i - \mu)^T (x_j - \mu) = (x_i - \mu)^T (x_j - \mu)$$

$$(x_i - \mu)^T (x_j - \mu) = (x_i - \mu)^T (x_j - \mu)$$

$$(x_i - \mu)^T (x_j - \mu) = (x_i - \mu)^T (x_j - \mu)$$

$$(x_i - \mu)^T (x_j - \mu) = (x_i - \mu)^T (x_j - \mu)$$

$$(x_i - \mu)^T (x_j - \mu) = (x_i - \mu)^T (x_j - \mu)$$

$$(x_i - \mu)^T (x_j - \mu) = (x_i - \mu)^T (x_j - \mu)$$

$$(x_i - \mu)^T (x_j - \mu) = (x_i - \mu)^T (x_j - \mu)$$

$$(x_i - \mu)^T (x_j - \mu) = (x_i - \mu)^T (x_j - \mu)$$

$$(x_i - \mu)^T (x_j - \mu) = (x_i - \mu)^T (x_j - \mu)$$

$$(x_i - \mu)^T (x_j - \mu) = (x_i - \mu)^T (x_j - \mu)$$

$T(y) \in \mathbb{R}^{K-1}$  - one-hot vectors

$$T(1) = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, T(2) = \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}, T(K-1) = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix}, T(K) = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$(T(y))_i = i^{\text{th}}$  element of vector  $T(y)$

indicator function -  $1\{\cdot\} = 1$  if true, 0 if false  $\therefore (T(y))_i = 1\{y=i\}$

$$E[(T(y))_i] = P(y=i) = \phi_i$$

$$\begin{aligned} p(y; \phi) &= \phi_1^{1\{y=1\}} \phi_2^{1\{y=2\}} \dots \phi_K^{1\{y=K\}} \\ &= \phi_1^{(T(y))_1} \phi_2^{(T(y))_2} \dots \phi_K^{1-\sum_{i=1}^{K-1} (T(y))_i} \\ &= \exp((T(y))_1 \log \phi_1 + (T(y))_2 \log \phi_2 + \dots + (T(y))_{K-1} \log \phi_{K-1}) \\ &= \exp((T(y))_1 \log \frac{\phi_1}{\phi_K} + (T(y))_2 \log \frac{\phi_2}{\phi_K} + \dots + (T(y))_{K-1} \log \frac{\phi_{K-1}}{\phi_K} + \log \phi_K) \end{aligned}$$

$$\eta = \begin{bmatrix} \log \frac{\phi_1}{\phi_K} \\ \log \frac{\phi_2}{\phi_K} \\ \vdots \\ \log \frac{\phi_{K-1}}{\phi_K} \end{bmatrix}, a(\eta) = -\log \phi_K, b(y) = 1$$

link function (for  $i=1, \dots, k$ ) -  $\eta_i = \log \frac{\phi_i}{\phi_K}, \eta_K = \log \frac{\phi_K}{\phi_K} = 0$

response function -  $e^{\eta_i} = \frac{\phi_i}{\phi_K}$

substitution

$$\phi_K e^{\eta_i} = \phi_i$$

$$\phi_K \sum_{i=1}^K e^{\eta_i} = \sum_{i=1}^K \phi_i = 1 \therefore \phi_k = \frac{1}{\sum_{i=1}^K e^{\eta_i}}$$

$$\phi_i = \frac{e^{\eta_i}}{\sum_{j=1}^K e^{\eta_j}} \quad \text{softmax function}$$

$$p(y=i|x; \theta) = \phi_i = \frac{e^{\eta_i}}{\sum_{j=1}^K e^{\eta_j}} = \frac{e^{\theta_i^T x}}{\sum_{j=1}^K e^{\theta_j^T x}}$$

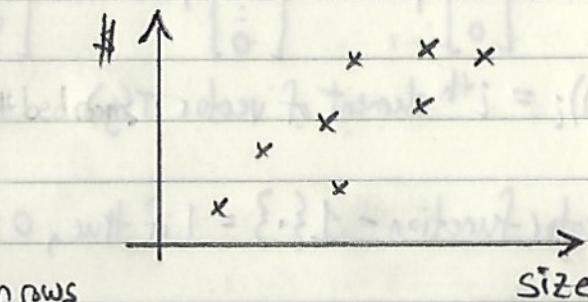
$$h_\theta(x) = E[T(y)|x; \theta] = E \left[ \begin{array}{c|c} 1\{y=1\} & \\ 1\{y=2\} & x^T \theta \\ \vdots & \\ 1\{y=K\} & \end{array} \right] = \begin{bmatrix} \phi_1 \\ \phi_2 \\ \vdots \\ \phi_{K-1} \end{bmatrix}$$

Supervised learning:  $x \rightarrow y$   
regressions vs. classification

Housing data = training set

size #bed rooms price ('000)

$x^{(1)}$	1	2104	4	\$400
$x^{(2)}$	1	1416	3	\$232
$\vdots$	$\vdots$	$\vdots$	$\vdots$	



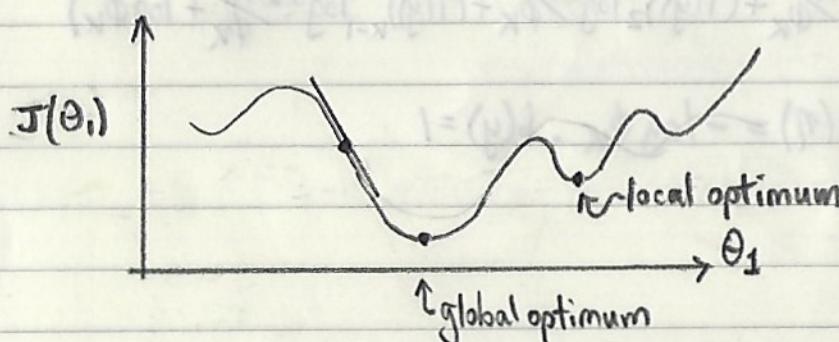
$$h(x) = \theta_0 + \theta_1 x \leftarrow \text{affine function}$$

$x_1$ : size

$x_2$ : # bedrooms

$x_1^{(2)} = 1416$

$$\vec{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix}, \vec{x} = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix}, x_0 = 1$$



## Machine Learning

CS229: L2

Linear Regression, Gradient Descent, Normal Equations

continuous value prediction - regression problem

$n$  - # training examples,  $(x, y)$   $d$  - # features

$x$  - input variables/features

$y$  - output variable/target

$(x^{(i)}, y^{(i)})$  -  $i^{\text{th}}$  training example

$\uparrow x, \theta$  are  $(d+1)$  dimensional

training examples

"hypothesis"

learning algorithm

linear

input  $\rightarrow h$  estimate

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 = \sum_{j=0}^2 \theta_j x_j = \theta^T x$$

$\theta$  - parameters

predicted value

$$\uparrow x_0 = 1 \\ = \sum_{j=0}^d \theta_j x_j$$

$$\min_{\theta} \frac{1}{2} \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)})^2 = \min_{\theta} J(\theta)$$

start with  $\theta = \vec{0}$  (init to zeros)

Keep changing  $\theta$  to reduce  $J(\theta)$

cost function

$$\begin{aligned} \text{gradient descent: } \theta_j &:= \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) := \theta_j - \alpha (h_{\theta}(x) - y) x_j \\ \frac{\partial}{\partial \theta_j} J(\theta) &= \frac{\partial}{\partial \theta_j} \frac{1}{2} (h_{\theta}(x) - y)^2 \\ &= 2 \frac{1}{2} (h_{\theta}(x) - y) \frac{\partial}{\partial \theta_j} (h_{\theta}(x) - y) \\ &= (h_{\theta}(x) - y) \frac{\partial}{\partial \theta_j} (\theta_0 x_0 + \dots + \theta_d x_d - y) \\ &= (h_{\theta}(x) - y) x_j \end{aligned} \quad \leftarrow \text{single training example}$$

repeat until convergence = batch gradient descent (works for small  $n$ )

$$\theta_j := \theta_j - \frac{\alpha}{n} \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_j \quad \leftarrow n \text{ training examples}$$

ordinary least squares - one global min.  
with linear  $h$

repeat until convergence: Stochastic gradient descent

for  $i = 1 \text{ to } n$ :

$$\theta_j := \theta_j - \alpha (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} \quad (\text{for all } j = 0 \dots d)$$

won't converge to global minimum - need to "time out"

Normal equations

$$\nabla_{\theta} J(\theta) = \begin{bmatrix} \frac{\partial J}{\partial \theta_0} \\ \vdots \\ \frac{\partial J}{\partial \theta_d} \end{bmatrix} \in \mathbb{R}^{d+1} \quad \text{gradient descent:} \\ \theta := \theta - \alpha \nabla_{\theta} J \\ \in \mathbb{R}^{d+1}$$

$$f: \mathbb{R}^{n \times d} \mapsto \mathbb{R}, f(A) \text{ where } A \in \mathbb{R}^{n \times d}$$

$$\nabla_A f(A) = \begin{bmatrix} \frac{\partial f}{\partial A_{11}} & \dots & \frac{\partial f}{\partial A_{1d}} \\ \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial A_{m1}} & \dots & \frac{\partial f}{\partial A_{md}} \end{bmatrix}, \text{tr } A = \sum_{i=1}^d A_{ii} - \text{sum diag. elements}$$

$$X = \begin{bmatrix} -(x^{(1)})^T \\ \vdots \\ -(x^{(n)})^T \end{bmatrix}, X\theta = \begin{bmatrix} -x^{(1)\top}\theta \\ \vdots \\ -x^{(n)\top}\theta \end{bmatrix} = \begin{bmatrix} h_\theta(x^{(1)}) \\ \vdots \\ h_\theta(x^{(n)}) \end{bmatrix}$$

$$\vec{y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix}, X\theta - \vec{y} = \begin{bmatrix} h_\theta(x^{(1)}) - y^{(1)} \\ \vdots \\ h_\theta(x^{(n)}) - y^{(n)} \end{bmatrix} \uparrow \text{recall } \vec{z}^T \vec{z} = \sum_i z_i^2 \downarrow \text{vector}$$

$$\nabla_{\theta} \text{tr} \theta \theta^T X^T X = \nabla_{\theta} \text{tr} \theta I \theta^T X^T X = X^T X \theta I + X^T X \theta I \\ \stackrel{\substack{\uparrow \uparrow \uparrow \\ A \quad B \quad C}}{=} X^T X \theta + X^T X \theta$$

$$\nabla_{\theta} \text{tr} y^T X \theta = X^T y$$

$$X^T X \theta = X^T \vec{y}$$

$$\theta = (X^T X)^{-1} X^T \vec{y}$$

not invertible if dependent features -  
pseudo inverse

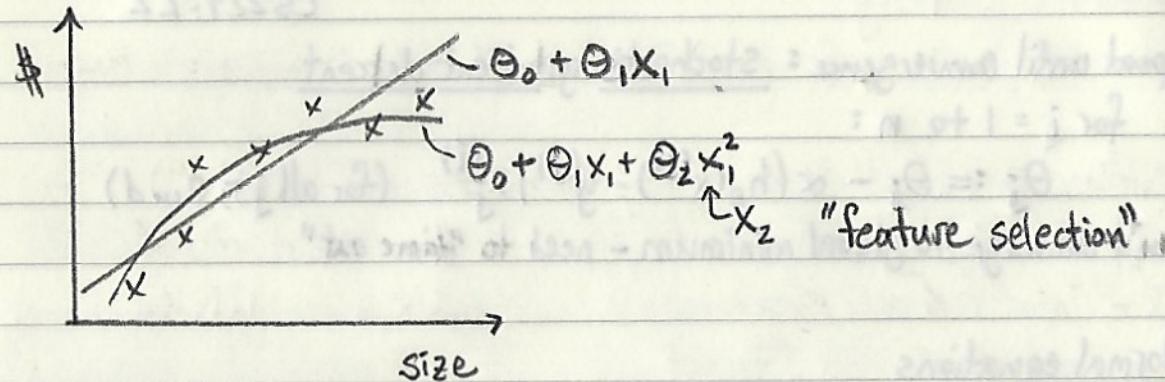
$$\frac{1}{2} (X\theta - \vec{y})^T (X\theta - \vec{y}) = \frac{1}{2} \sum_{i=1}^n (h_\theta(x^{(i)}) - y^{(i)})^2 = J(\theta)$$

$$\begin{aligned} \nabla_{\theta} J(\theta) &= \vec{0} = \nabla_{\theta} \frac{1}{2} (X\theta - \vec{y})^T (X\theta - \vec{y}) \\ &= \frac{1}{2} \nabla_{\theta} \text{tr} (\theta^T X^T X \theta - \theta^T X^T \vec{y} - \vec{y}^T X \theta + \vec{y}^T \vec{y}) \\ &= \frac{1}{2} [\nabla_{\theta} \text{tr} \theta \theta^T X^T X - \nabla_{\theta} \text{tr} \vec{y}^T X \theta - \nabla_{\theta} \text{tr} \vec{y}^T \vec{y}] \\ &= \frac{1}{2} [X^T X \theta + X^T X \theta - X^T \vec{y} - X^T \vec{y}] = X^T X \theta - X^T \vec{y} = 0 \end{aligned}$$

does not depend  
on  $\theta$

## Machine Learning

### Locally Weighted Regression, Probabilistic Interpretation, Logistic Regression



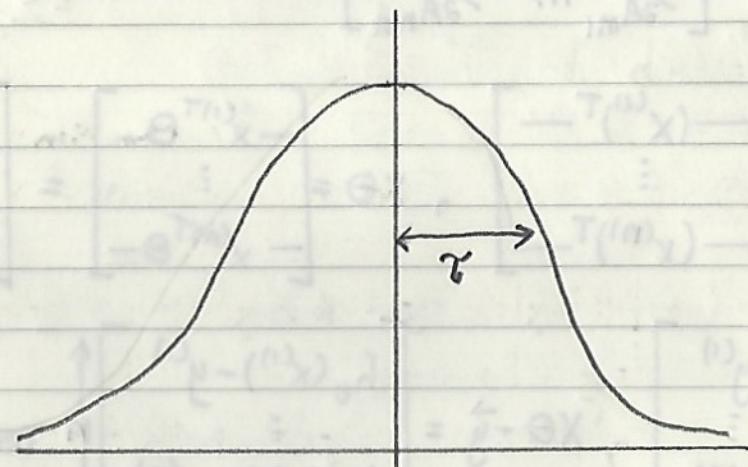
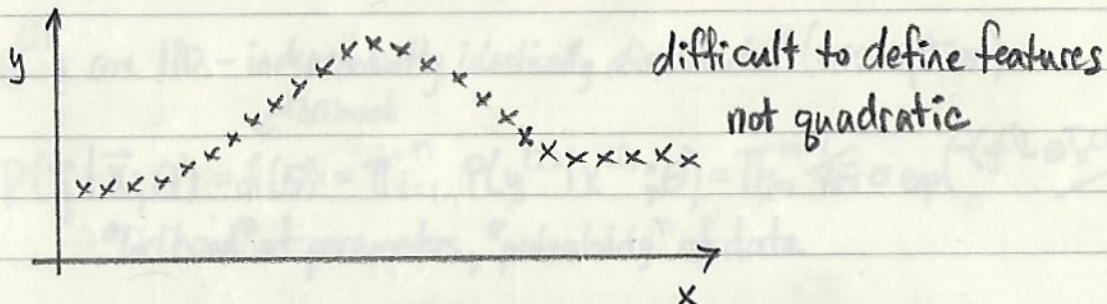
$$h_{\theta}(x) = \sum_{j=0}^d \theta_j x_j = \theta^T x \quad \text{where } x_0 = 1 \quad (x^{(i)}, y^{(i)}) \text{ i}^{\text{th}} \text{ example}$$

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)})^2 \quad \text{- quadratic cost} \quad x^{(i)} \in \mathbb{R}^{d+1}, x_0 = 1$$

$$\theta = (X^T X)^{-1} X^T y \quad \text{- closed form}$$

"Parametric" learning algorithm -  $\theta$ s fixed set of parameters fit to data  
↑ linear regression

"Non-parametric" learning algorithm - # of parameters grows (linearly) with n  
↑ locally weighted regression (loess)



to evaluate  $h$  at a certain  $x$

LR: fit  $\theta$  to min.  $\sum_i (y^{(i)} - \theta^T x^{(i)})^2$  then return  $\theta^T x$

LWR = apply LR to data points only in vicinity of  $x$

fit  $\theta$  to min.  $\sum_i w^{(i)} (y^{(i)} - \theta^T x^{(i)})^2$  where  $w^{(i)} = \exp\left(-\frac{(x^{(i)} - x)^2}{2\gamma^2}\right)$

if  $|x^{(i)} - x|$  is small then  $w^{(i)} \approx 1$  ↑ other weight algorithms

if  $|x^{(i)} - x|$  is large then  $w^{(i)} \approx 0$  are possible

$\gamma$  is "bandwidth" - controls width of bell curve

need to run new fitting for every prediction - not "building model"

can still overfit or underfit

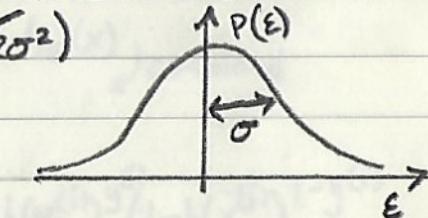
↑ expensive

$w^{(i)}$  = weighting function

Probabalistic Interpretation of Least Squares

assume  $y^{(i)} = \theta^T x^{(i)} + \varepsilon^{(i)}$ ,  $\varepsilon^{(i)}$  - error, unmodeled effects + noise  
 $\varepsilon^{(i)} \sim N(0, \sigma^2)$  - Gaussian (central limit theorem)

$$\text{density } P(\varepsilon^{(i)}) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(\varepsilon^{(i)})^2}{2\sigma^2}\right)$$



$$P(y^{(i)} | x^{(i)}; \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)$$

not random variable

$$y^{(i)} | x^{(i)}; \theta \sim N(\theta^T x, \sigma^2)$$

$\varepsilon^{(i)}$ 's are IID - independently identically distributed (assumption)  
 ↘ likelihood

$$P(\vec{y} | \vec{x}; \theta) = \mathcal{L}(\theta) = \prod_{i=1}^n P(y^{(i)} | x^{(i)}; \theta) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)$$

"likelihood" of parameters, "probability" of data

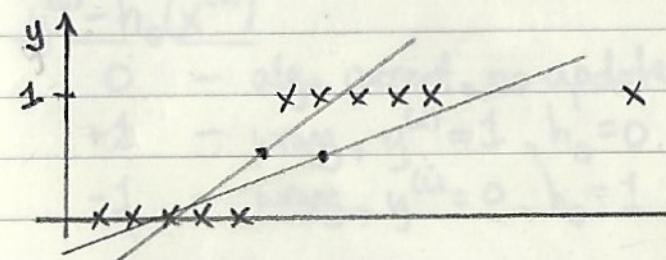
maximum likelihood: choose  $\theta$  to max.  $\mathcal{L}(\theta) = P(\vec{y} | \vec{x}; \theta)$

$$\begin{aligned} \log \text{likelihood } \ell(\theta) &= \log \mathcal{L}(\theta) = \log \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right) \\ &= \sum_{i=1}^n \log \left[ \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right) \right] \\ &= n \log \frac{1}{\sqrt{2\pi}\sigma} - \sum_{i=1}^n \frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2} \end{aligned}$$

$$\text{max. } \ell(\theta) \text{ same as min. } \frac{1}{2} \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2 = J(\theta)$$

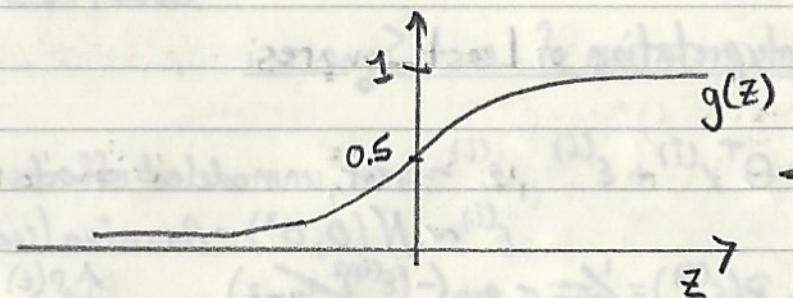
↑ independent of  $\sigma^2$

Classification  $y \in \{0, 1\}$  - binary



LR generally does not work, although it can sometimes

CS 229:L3

Logistic Regression

$y \in \{0, 1\}$ ,  $h_{\theta}(x) \in [0, 1]$  ← bound predictions to between 0 and 1  
choose:  $h_{\theta}(x) = g(\theta^T x) = \frac{1}{1+e^{-\theta^T x}}$   
↑ sigmoid,  $g(z) = \frac{1}{1+e^{-z}}$  - logistic function

$$P(y=1|x; \theta) = h_{\theta}(x), P(y=0|x; \theta) = 1 - h_{\theta}(x) \quad \text{combined}$$

$$P(y|x; \theta) = h_{\theta}(x)^y (1-h_{\theta}(x))^{1-y}$$

$$\mathcal{L}(\theta) = P(\bar{y}|x; \theta) = \prod_{i=1}^n p(y^{(i)}|x^{(i)}; \theta) = \prod_{i=1}^n h(x^{(i)})^{y^{(i)}} (1-h(x^{(i)}))^{1-y^{(i)}}$$

$$l(\theta) = \log \mathcal{L}(\theta) = \sum_{i=1}^n y^{(i)} \log h(x^{(i)}) + (1-y^{(i)}) \log (1-h(x^{(i)}))$$

gradient ascent -  $\theta := \theta + \alpha \nabla_{\theta} l(\theta)$

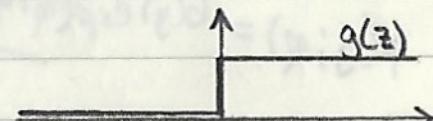
$$\frac{\partial}{\partial \theta_j} l(\theta) = \sum_{i=1}^n (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$$

$$\theta_j := \theta_j + \alpha \frac{\partial}{\partial \theta_j} l(\theta)$$

$$:= \theta_j + \alpha \sum_{i=1}^n (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$$

Perceptron

$$g(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{otherwise} \end{cases} \quad \text{- step function}$$

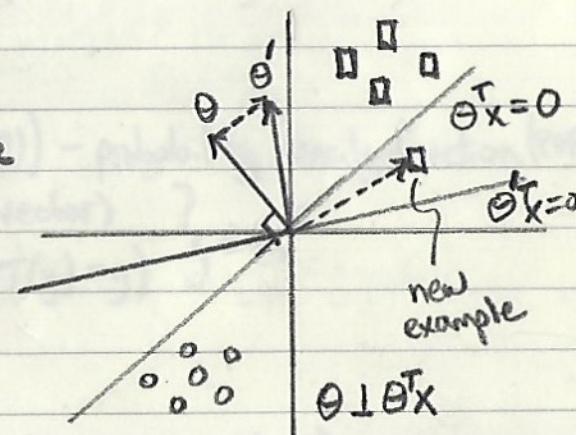


$$h(x) = g(\theta^T x) = \begin{cases} 1 & \text{if } \theta^T x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

$$\theta_j := \theta_j + \alpha (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$$

$$y^{(i)} - h_{\theta}(x^{(i)})$$

- 0 - algo correct, no update
- +1 - wrong,  $y^{(i)} = 1, h_{\theta} = 0$
- 1 - wrong,  $y^{(i)} = 0, h_{\theta} = 1$



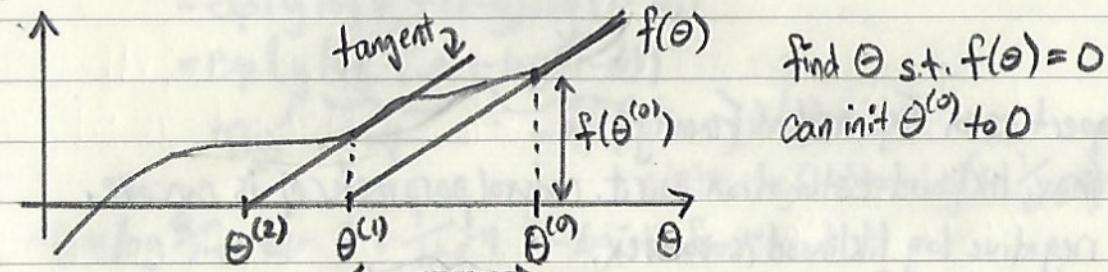
## Machine Learning

### Newton's Method, Exponential Families, Generalized Linear Models

CS229-L4

quadratic convergence - very fast

Newton's method runs faster than gradient ascent for logistic regression.



$$\theta^{(t+1)} := \theta^{(t)} - \frac{f(\theta^{(t)})}{f'(\theta^{(t)})}, \Delta = \frac{f(\theta^{(t)})}{f'(\theta^{(t)})}, \theta^{(1)} = \theta^{(0)} - \frac{f(\theta^{(0)})}{f'(\theta^{(0)})}$$

$$\theta^{(t+1)} := \theta^{(t)} - \frac{f(\theta^{(t)})}{f'(\theta^{(t)})} - \text{iterate}$$

$$l(\theta) \text{ want } \theta \text{ s.t. } l'(\theta) = 0 \therefore \theta^{(t+1)} = \theta^{(t)} - \frac{l'(\theta^{(t)})}{l''(\theta^{(t)})}$$

$$\theta^{(t+1)} := \theta^{(t)} - H^{-1} \nabla_{\theta} l(\theta) \quad - \theta \text{ is a vector}$$

$$\text{Hessian} - H_{ij} = \frac{\partial^2 l(\theta)}{\partial \theta_i \partial \theta_j} \in \mathbb{R}^{(d+1) \times (d+1)}$$

$H^{-1}$  expensive to compute with large d

$$p(y|x;\theta)$$

$y \in \mathbb{R}$  = Gaussian  $\rightarrow$  least squares

$y \in \{0, 1\}$  = Bernoulli  $\rightarrow$  logistic regression

Bernoulli( $\phi$ ) -  $P(y=1; \phi) = \phi \quad \} \text{Exponential Family Distributions}$

$$\mathcal{N}(\mu, \sigma^2)$$

Exponential Family

$p(y; \eta) = b(y) \exp(\eta^T T(y) - a(\eta))$  - probability density function (PDF)

$\eta$  - natural parameter (scalar/vector)  $\} \in \mathbb{R}^k$

$T(y)$  - sufficient statistic (usually  $T(y) = y$ )  $\} \in \mathbb{R}^k$

$y$  - data (scalar)

$b(y)$  - base measure (scalar)

$a(\eta)$  - log-partition function (makes PDF sum to 1)

$$\begin{aligned}-\log(1 - \frac{e^{-\eta}}{1+e^{-\eta}}) &= -\log(\frac{e^{-\eta}}{1+e^{-\eta}}) \\ &= -\log(\frac{1}{1+e^{\eta}}) = \log(1+e^{\eta})\end{aligned}$$

Properties of Exponential Family:

- a) max. likelihood estimation w.r.t. natural parameter  $\eta$  is concave.  
negative log likelihood is convex.
- b)  $E[y; \eta] = \frac{\partial}{\partial \eta} a(\eta)$
- c)  $\text{Var}[y; \eta] = \frac{\partial^2}{\partial \eta^2} a(\eta)$

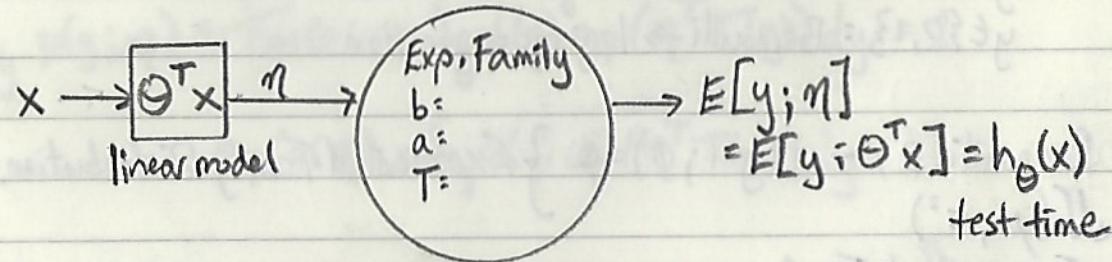
Real : Gaussian

Binary : Bernoulli

Count : Poisson

$R^+$  : Gamma, Exponential

Dist<sup>n</sup> : Beta, Dirichlet



$$\text{training} = \max_{\Theta} \log p(y^{(i)}; \Theta^T x^{(i)})$$

$$\text{learning: } \Theta_j := \Theta_j + \alpha (y^{(i)} - h_\Theta(x^{(i)})) x_j^{(i)}$$

$\uparrow$  plug in appropriate  $h_\Theta(x)$

$$\text{Bernoulli}(\phi) - p(y=1; \phi) = \phi, p(y=0; \phi) = 1-\phi$$

$$P(y; \phi) = \phi^y (1-\phi)^{1-y}$$

=  $\exp(\log(\phi^y (1-\phi)^{1-y}))$  - exp and log cancel each other

$$= \exp(y \log \phi + (1-y) \log(1-\phi))$$

$$= \exp(y \underbrace{\log \frac{\phi}{1-\phi}}_{\eta} + \underbrace{\log(1-\phi)}_{-a(\eta)})$$

$$T(y) \quad \eta \quad -a(\eta) \quad b(y) = 1$$

$$\eta = \log \frac{\phi}{1-\phi} \Rightarrow \phi = \frac{1}{1+e^\eta} \leftarrow \text{logistic function}$$

Gaussian:  $N(\mu, \sigma^2)$ , set  $\sigma^2=1$  since no impact on max. likelihood estimate

$$\frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(y-\mu)^2\right) = p(y|\mu) = \underbrace{\frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}y^2\right)}_{b(y)} \exp\left(\mu y - \frac{1}{2}\mu^2\right) \underbrace{\eta \uparrow \underbrace{T(y)}_{\text{link function}}}_{-a(\eta)} \underbrace{a(\eta) = \frac{1}{2}\mu^2}_{\text{variance}}$$

### Generalized Linear Model (GLM)

assume:

$$1) y|x; \Theta \sim \text{ExpFamily}(\eta)$$

$$2) \text{given } x, \text{ goal is to output } E[y|x; \Theta] = h_\Theta(x)$$

$$\text{want } h(x) = E[T(y)|x]$$

$$3) \eta = \Theta^T x - \text{linear relationship between } \eta \text{ and } x \text{ (design choice)}$$

$\hookrightarrow \eta_i = \Theta_i^T x$  if  $\eta \in \mathbb{R}^K$        $\Theta \in \mathbb{R}^d, x \in \mathbb{R}^d$

$$\text{Bernoulli} = y|x; \Theta \sim \text{ExpFamily}(\eta)$$

$$\text{for fixed } x, \Theta, \text{ algorithm output } h_\Theta(x) = E[y|x; \Theta] = p(y=1|x; \Theta) = \phi = \frac{1}{1+e^{-\eta}} = \frac{1}{1+e^{-\Theta^T x}}$$

$$g(\eta) = E[y; \eta] = \frac{1}{1+e^{-\eta}} - \text{canonical response function}$$

$$g^{-1} - \text{canonical link function}$$

$$\mu = E[y; \eta] = g(\eta) = \frac{\partial}{\partial \eta} a(\eta)$$

$$\eta = g^{-1}(\mu)$$

3 parameters:

- 1) model
- 2) natural
- 3) Canonical

$$\Theta \xrightarrow{\Theta^T X} \eta \xrightarrow{g} \begin{array}{l} \phi = \text{Bernoulli} \\ \mu = \text{Gaussian} \\ \lambda = \text{Poisson} \end{array}$$

logistic regression:

$$h_\theta(x) = E[y|x; \theta] = \phi = \frac{1}{1+e^{-\eta}} = \frac{1}{1+e^{-\Theta^T x}}$$

Multinomial:  $y \in \{1, \dots, K\}$

parameters:  $\phi_1, \phi_2, \dots, \phi_{K-1}, \phi_K = 1 - (\phi_1 + \dots + \phi_{K-1})$

$$P(y=i) = \phi_i \quad \text{↑ redundant parameter}$$

$$T(1) = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, T(2) = \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \dots \in \mathbb{R}^{K-1}, T(K-1) = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix}, T(K) = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$1\{y \text{ True}\} = 1, 1\{y \text{ False}\} = 0$  - indicator function

$$T(y)_i = 1\{y=i\}$$

$$P(y) = \phi_1^{1\{y=1\}} \phi_2^{1\{y=2\}} \dots \phi_K^{1\{y=K\}} = \phi_1^{T(y)_1} \phi_2^{T(y)_2} \dots \phi_K^{1 - \sum_{j=1}^{K-1} T(y)_j}$$

$$= b(y) \exp(\eta^T T(y) - a(\eta))$$

$$\eta = \begin{bmatrix} \log(\phi_1/\phi_K) \\ \vdots \\ \log(\phi_{K-1}/\phi_K) \end{bmatrix} \in \mathbb{R}^{K-1}, a(\eta) = -\log(\phi_K), b(y) = 1$$

$$\phi_i = \frac{e^{\eta_i}}{1 + \sum_{j=1}^{K-1} e^{\eta_j}} \quad (i = 1, \dots, K-1)$$

$$= \frac{e^{\Theta_i^T x}}{1 + \sum_{j=1}^{K-1} e^{\Theta_j^T x}} \quad (\text{using } \eta_i = \Theta_i^T x)$$

$$h_\theta(x) = E[T(y)|x; \theta] = E\left[\begin{array}{c} 1\{y=1\} \\ \vdots \\ 1\{y=K-1\} \end{array} \middle| x; \theta\right] = \begin{bmatrix} \phi_1 \\ \vdots \\ \phi_{K-1} \end{bmatrix}$$

$$= \begin{bmatrix} \frac{e^{\Theta^T x}}{1 + \sum_{j=1}^{K-1} e^{\Theta_j^T x}} \\ \vdots \\ \frac{e^{\Theta_{K-1}^T x}}{1 + \sum_{j=1}^{K-1} e^{\Theta_j^T x}} \end{bmatrix}$$

Softmax regression (multiclass classification)

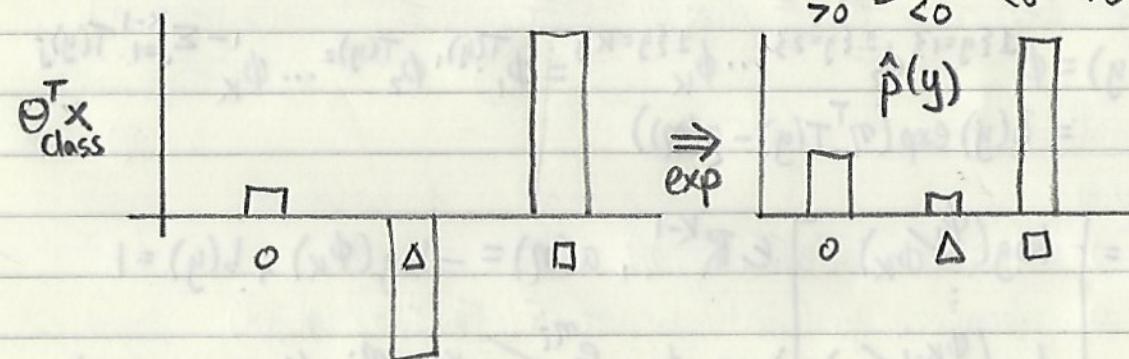
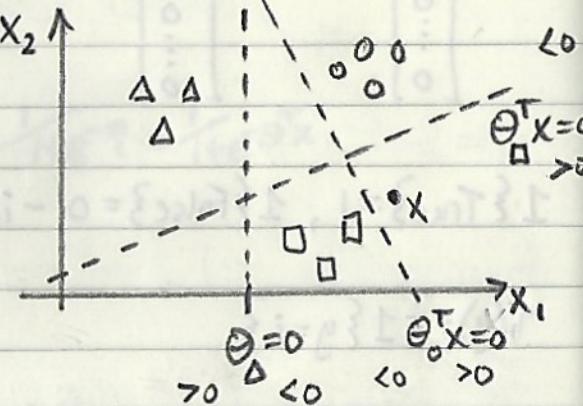
Cross Entropy Minimization

$$x^{(i)} \in \mathbb{R}^d$$

$$y \in \{0, 1\}^K \text{ eg. } (0, 0, 1, 0) - \text{one-hot vector}$$

$$\theta_{\text{class}} \in \mathbb{R}^d \quad (K \text{ such})$$

$$\begin{bmatrix} -\theta_1 & - \\ \vdots & \\ -\theta_K & - \end{bmatrix}$$



$$\text{cross entropy } (p, \hat{p}) = \sum_{y \in \{0, \Delta, 1\}} p(y) \log \hat{p}(y)$$

$$= -\log \hat{p}(y_0)$$

$$= -\log \left( \frac{e^{\theta_0^T x}}{\sum_{\text{class} \in \{0, \Delta, 1\}} e^{\theta_{\text{class}}^T x}} \right)$$

Andrew Ng  
Moses Charikar  
(F19) 10/10 produced without AI

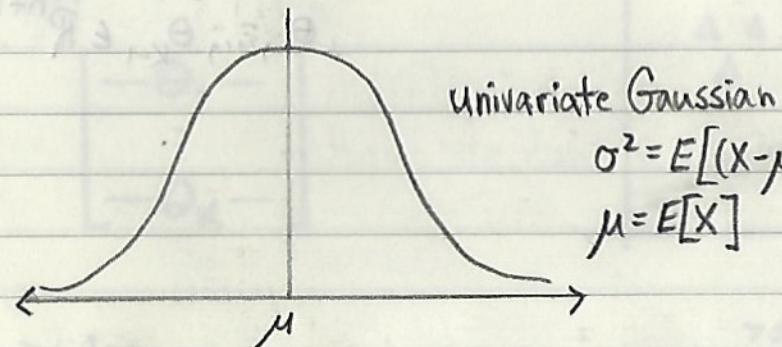
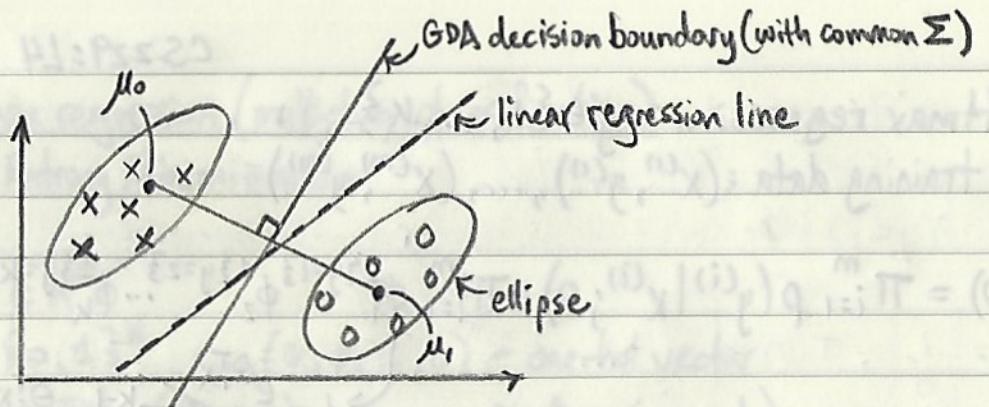
CS 229: L4

Softmax regression =  $y \in \{1, \dots, K\}$   
training data:  $(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$

$$L(\theta) = \prod_{i=1}^m p(y^{(i)} | x^{(i)}; \theta) = \prod_{i=1}^m \phi_1^{1 \leq y^{(i)} \leq 3} \phi_2^{1 \leq y^{(i)} \leq 2} \cdots \phi_K^{1 \leq y^{(i)} \leq K}$$

$$\phi_j = \frac{e^{\theta_j^T x}}{1 + \sum_{j=1}^{K-1} e^{\theta_j^T x}}$$

$$\theta_1, \dots, \theta_{K-1} \in \mathbb{R}^{n+1}$$



$$\sigma^2 = E[(X-\mu)^2] = E[X^2] - \mu^2$$

$$\mu = E[X]$$

$$E[Z] = \bar{\mu} \quad \leftarrow$$

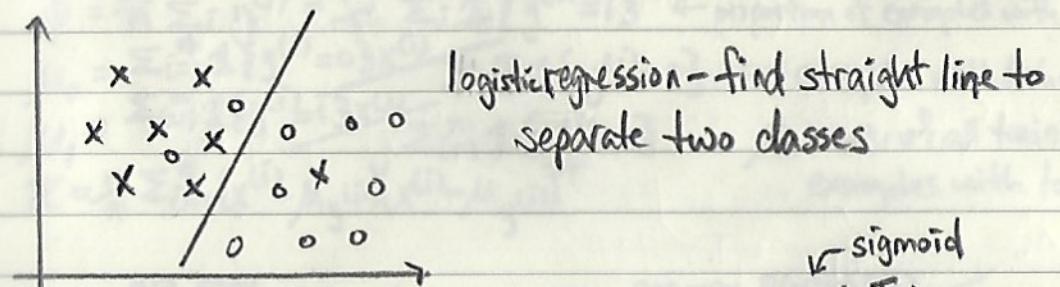
$$\Sigma = E[ZZ^T] - (E[Z])(E[Z])^T$$

parameters:  $\mu_0, \mu_1, \Sigma, \phi$   
 $\mu_0, \mu_1 \in \mathbb{R}^d$   
 $\Sigma \in \mathbb{R}^{d \times d}$   
 $\phi \in [0, 1]$

$$\mathcal{L}(\phi, \mu_0, \mu_1, \Sigma) = \prod_{i=1}^n p(x^{(i)}, y^{(i)}; \phi, \mu_0, \mu_1, \Sigma)$$

$$= \prod_{i=1}^n p(x^{(i)}|y^{(i)}; \mu_0, \Sigma) p(y^{(i)})$$

$$\mathcal{L}(\theta) = \prod_{i=1}^n p(y^{(i)}|x^{(i)}; \theta) \leftarrow \text{log. reg.}$$



discriminative - learns  $p(y|x)$  or learns  $h_\theta(x) \in \{0, 1\}$  directly

generative -  $p(x|y)$  ( $\neq p(y)$  - class prior)

features  $\uparrow$  class label

$$\text{using Bayes' Rule: } p(y=1|x) = \frac{p(x|y=1)p(y=1)}{p(x)}$$

$$p(x) = p(x|y=1)p(y=1) + p(x|y=0)p(y=0)$$

↑ not necessary since  $\arg \max_y p(y|x) = \arg \max_y p(x|y)p(y)$

assume  $X \in \mathbb{R}^d$ , continuous valued

Gaussian Discriminant Analysis =  $p(x|y)$  is Gaussian  $\leftarrow$  assume

multivariate Gaussian:  $\Sigma \sim N(\mu, \Sigma)$ ,  $\mu \in \mathbb{R}^d$ ,  $\Sigma \in \mathbb{R}^{d \times d}$ ,  $Z \in \mathbb{R}^d$

$$p(x) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)\right)$$

mean  $\uparrow$  Covariance  $\Sigma = E[(Z-\mu)(Z-\mu)^T]$   
matrix  $\Sigma_{ij} = E[(z_i - \mu_i)(z_j - \mu_j)]$

$$p(y) = \phi^y (1-\phi)^{1-y} \sim \text{Bernoulli}(\phi) \quad p(y=1) = \phi$$

$$p(x|y=0) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (x-\mu_0)^T \Sigma^{-1} (x-\mu_0)\right) \sim N(\mu_0, \Sigma)$$

$$p(x|y=1) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (x-\mu_1)^T \Sigma^{-1} (x-\mu_1)\right) \sim N(\mu_1, \Sigma)$$

$\leftarrow$  joint likelihood

$$l(\phi, \mu_0, \mu_1, \Sigma) = \log \prod_{i=1}^n p(x^{(i)}, y^{(i)}; \phi, \mu_0, \mu_1, \Sigma)$$

$$= \log \prod_{i=1}^n p(x^{(i)}|y^{(i)}; \mu_0, \mu_1, \Sigma) p(y^{(i)}; \phi)$$

$$l(\theta) = \log \prod_{i=1}^n p(y^{(i)}|x^{(i)}; \theta) \leftarrow \text{conditional likelihood, logistic regression}$$

$\uparrow$  for comparison

maximize  $\ell$  w.r.t.  $\phi, \mu_0, \mu_1, \Sigma$ : no iteration

$$\phi = \frac{1}{n} \sum_{i=1}^n y^{(i)} = \frac{1}{n} \sum_i \mathbb{1}\{y^{(i)}=1\} \leftarrow \text{proportion of examples with } y^{(i)}=1$$

~~$$\mu_0 = \frac{\sum_{i=1}^n \mathbb{1}\{y^{(i)}=0\}}{\sum_{i=1}^n \mathbb{1}\{y^{(i)}=0\}}$$~~ 
$$\leftarrow \# \text{examples with label } 0$$

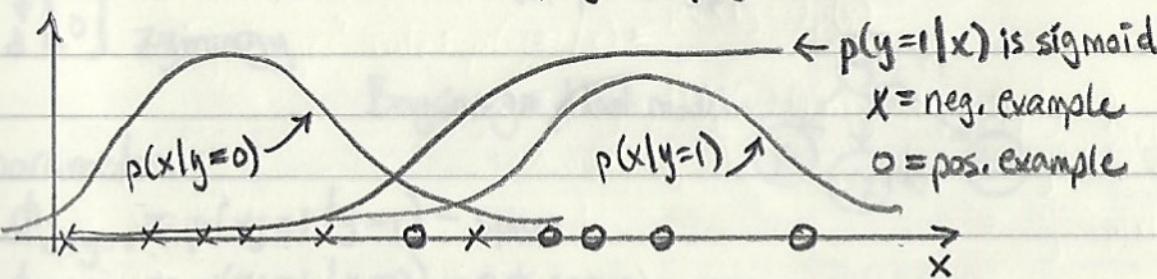
~~$$\mu_1 = \frac{\sum_{i=1}^n \mathbb{1}\{y^{(i)}=1\}}{\sum_{i=1}^n \mathbb{1}\{y^{(i)}=1\}}$$~~ 
$$\leftarrow \mu_0 = \text{avg. of all training examples with label } 0$$

~~$$\Sigma = \frac{1}{n} \sum_{i=1}^n (x^{(i)} - \mu_{y^{(i)}})(x^{(i)} - \mu_{y^{(i)}})^T$$~~

$$\text{predict: } \arg \max_y P(y|x; \phi, \mu_0, \mu_1, \Sigma) = \arg \max_y \frac{P(x|y) P(y)}{P(x)}$$

$$\text{if } P(y) \text{ is uniform} = \arg \max_y P(x|y)$$

$$\leftarrow p(y=0) = p(y=1)$$



$$p(y=1|x) = \frac{p(x|y=1)p(y=1)}{p(x)} \rightarrow p(y=1) = \phi$$

$$\rightarrow p(x) = p(x|y=0)p(y=0) + p(x|y=1)p(y=1)$$

↙ stronger assumption

$x|y \sim \text{Gaussian}$  ↗ logistic posterior for  $p(y=1|x)$

$$\begin{aligned} x|y=1 &\sim \text{Poisson}(\lambda_1) \\ x|y=0 &\sim \text{Poisson}(\lambda_0) \end{aligned} \Rightarrow p(y=1|x) \text{ is logistic} \Leftrightarrow \begin{cases} x|y=1 \sim \text{Exp Family}(\eta_1) \\ x|y=0 \sim \text{Exp family}(\eta_0) \end{cases}$$

Logistic Regression assumes

$$p(y=1|x) = \frac{1}{1+e^{-\Theta^T x}} \quad (x_0=1)$$

Gaussian Discriminate Analysis makes a stronger assumption (e.g.  $x|y \sim \text{Gaussian}$ ) than logistic regression and so can work with less training data.  
Logistic Regression very robust for different distributions

$$p(y=1|x; \phi, \mu_0, \mu_1, \Sigma) = \frac{1}{1+exp(-\Theta^T x)}, \Theta \text{ approx. function of } \phi, \mu_0, \mu_1, \Sigma$$

Naive Bayes:  $y \in \{0, 1\} \leftarrow$  spam email?

↙ vocabulary

$X =$	$\begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \\ \vdots \\ 0 \end{bmatrix}$	$a \in \{0, 1\}^d, d = 50,000+, 2^{50000}$ possible values of $x$
	$aardvark$	$x_i \in \{0, 1\}$
	$aardwolf$	naive assumption: $x_i$ 's are conditionally independent given $y$
	$\vdots$	extremely strong
	$buy$	independent given $y$
	$\vdots$	$p(x_1, x_2, \dots, x_{50000}   y) = p(x_1   y)p(x_2   y, x_1)p(x_3   y, x_1, x_2) \dots$
	$CS229$	$= p(x_1   y)p(x_2   y)p(x_3   y) \dots p(x_{50000}   y)$
	$\vdots$	$= \prod_{i=1}^d p(x_i   y)$
	$zymurgy$	

parameters =

$$\phi_{j|y=1} = p(x_j = 1 | y = 1) - \text{spam}$$

$$\phi_{j|y=0} = p(x_j = 1 | y = 0) - \text{not spam}$$

$$\phi_y = p(y = 1) \quad \text{training set } \{(x^{(i)}, y^{(i)}) ; i = 1, \dots, m\}$$

joint likelihood  $\mathcal{L}(\phi_y, \phi_{j|y=1}, \phi_{j|y=0}) = \prod_{i=1}^n p(x^{(i)}, y^{(i)})$

$$\phi_{j|y=1} = \frac{\sum_{i=1}^n \mathbb{1}\{x_j^{(i)} = 1, y^{(i)} = 1\} + 1}{\sum_{i=1}^n \mathbb{1}\{y^{(i)} = 1\} + 2} \quad \begin{array}{l} \text{Laplace smoothing} \\ \text{fraction of spam emails where word } j \text{ appears} \end{array}$$

$$\phi_y = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{y^{(i)} = 1\} + 2 \quad \begin{array}{l} \text{because binary=spam or not spam} \end{array}$$

$$\rightarrow p(y=1|x) = \frac{p(x|y=1)p(y=1)}{p(x|y=1)p(y=1) + p(x|y=0)p(y=0)}$$

$$p(x|y=1) = \prod_{i=1}^{50000} p(x_i | y=1) \quad \begin{array}{l} \leftarrow \text{zero if any word unseen in } y=1 \\ (\text{e.g. } p(x_{30000} | y=1) = 0) \end{array}$$

Laplace smoothing =

$$\text{if } y \in \{1, \dots, K\}, P(y=j) = \frac{\sum_{i=1}^n \mathbb{1}\{y^{(i)} = j\} + 1}{n + K} \quad \begin{array}{l} \phi_{30000|y=1} \\ \uparrow \text{possible values of } y \end{array}$$

size	<400	400-800	800-1200	1200+
x	1	2	3	4

$$x \in \begin{bmatrix} 3500 \\ 100 \\ 5000 \\ 9765 \end{bmatrix} \in \mathbb{R}^{d_i} \quad d_i = \text{length of email}$$

$x_j \in \{1, \dots, 50000\}$

↑ very short email

order of words does not matter (no  $y^{(i)}$ )

$$p(x|y) \stackrel{\text{assume}}{=} \prod_{j=1}^{d_i} p(x_j|y)$$

parameters:

 $\phi_y = P(y=1)$ 
 $\phi_{K|y=0} = P(x_j=K|y=0)$ 
 $\phi_{K|y=1} = P(x_j=K|y=1)$

MLE:

$$\phi_{K|y=0} = \frac{\sum_{i=1}^n (\underbrace{\mathbb{I}\{y^{(i)}=0\}}_{\text{non-spam}} \sum_{j=1}^{d_i} \underbrace{\mathbb{I}\{x_j^{(i)}=K\}}_{\text{\# times word K appears}}) + 1}{\sum_{i=1}^n \mathbb{I}\{y^{(i)}=0\} \cdot d_i + 50000}$$

↑ size of dictionary

$$X = \begin{bmatrix} 1 & a \\ 0 & \text{aardvark} \\ 0 & \text{aardwolf} \\ \vdots & \vdots \\ 1 & 0 \end{bmatrix} \quad x_i \in \{0, 1\}$$

↑ vocabulary size      ↑ multivariate Bernoulli event model

Generative learning algo:  $P(x|y) = \prod_{i=1}^n P(x_i|y) \}_{\text{model}}$

$$\arg \max_y P(y|x) = \arg \max_y P(x|y) P(y)$$

if  $x_i \in \{1, \dots, K\}$ ,  $P(x|y) = \prod_{i=1}^n P(x_i|y)$

↑ discretize continuous variable      ↑ multinomial

multinomial event model: some words appear multiple times

$$(x_1^{(i)}, x_2^{(i)}, \dots, x_{d_i}^{(i)}) \text{ where } d_i = \# \text{ words in email} \leftarrow \text{different for every training example}$$

$$P(x, y) = (\prod_{i=1}^n P(x_i|y)) P(y), d \text{ is length of email}, y \in \{0, 1\}$$

parameters:  $\phi_{K|y=1} = P(x_j=K|y=1) = \sum_{i=1}^n \mathbb{I}\{y^{(i)}=1\} \sum_{j=1}^{d_i} \mathbb{I}\{x_j^{(i)}=K\} = \frac{1}{d_i} \sum_{j=1}^{d_i} \mathbb{I}\{x_j^{(i)}=K\}$

↑ smoothing

$$\phi_{K|y=0} = P(x_j=K|y=0)$$

↑ of all spam emails, fraction of words that were word K

$$\phi_y = P(y=1)$$

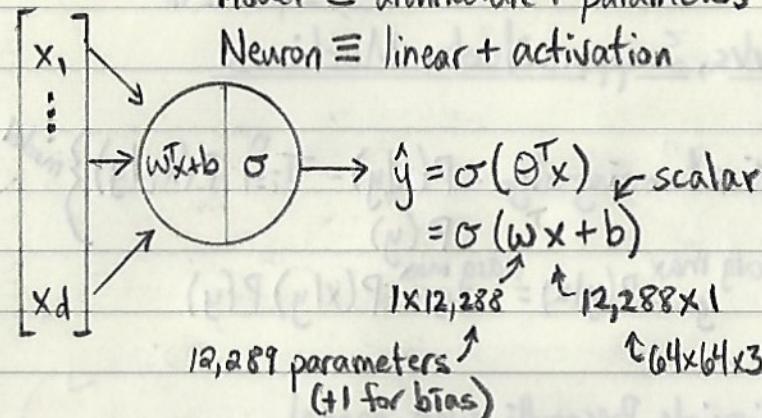
$$l(\phi_{K|y=1}, \phi_{K|y=0}, \phi_y) = \log \prod_{i=1}^n P(x^{(i)}, y^{(i)}; \phi_{K|y=1}, \phi_{K|y=0}, \phi_y)$$

$$= \log \prod_{i=1}^n \prod_{j=1}^{d_i} P(x_j^{(i)}, y^{(i)}; \phi_{K|y=1}, \phi_{K|y=0}) P(y^{(i)}; \phi_y)$$

performs better for text classification because takes into account number of times a word appears in a document

Model  $\equiv$  architecture + parameters

Neuron  $\equiv$  linear + activation



$$\hat{y} = \sigma(\theta^T x)$$

$$= \sigma(w x + b)$$

$$1 \times 12,288$$

scalar

$$12,288 \times 1$$

$$12,288 \times 1$$

$$64 \times 64 \times 3$$

$\Rightarrow$  scalar

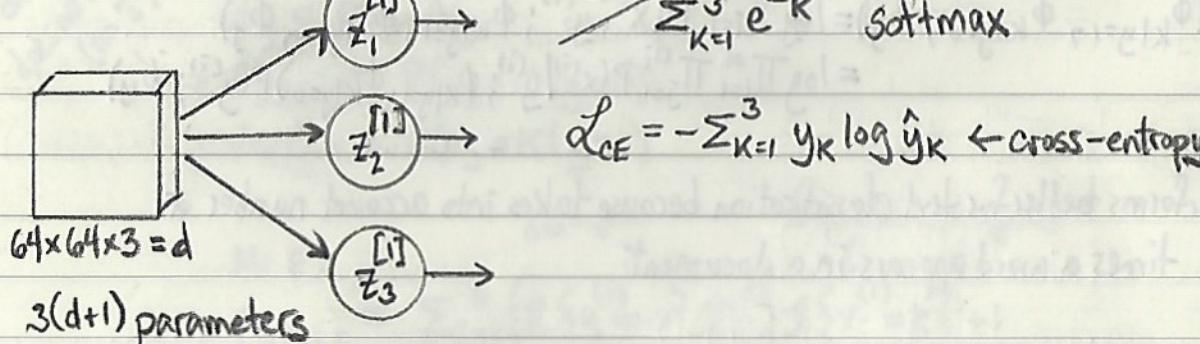
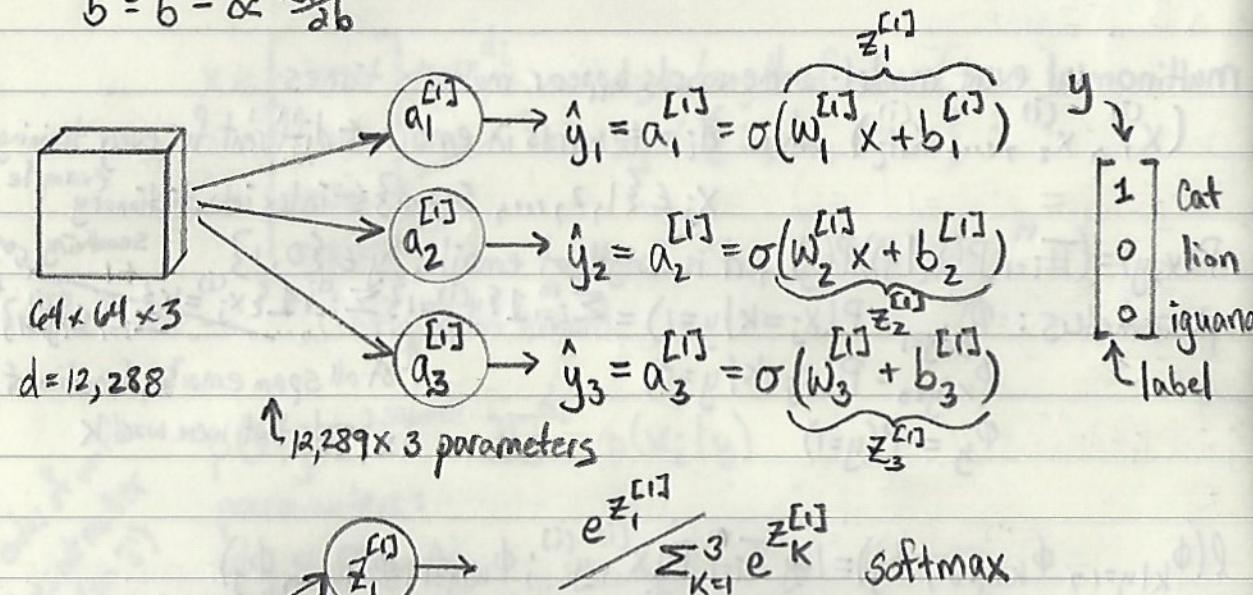
$\Rightarrow$  optimum  $w, b$

use  $\hat{y} = \sigma(w x + b)$  to predict

$$\mathcal{L} = -[y \log \hat{y} + (1-y) \log(1-\hat{y})] \leftarrow \text{log likelihood}$$

$$w = w - \alpha \frac{\partial \mathcal{L}}{\partial w}$$

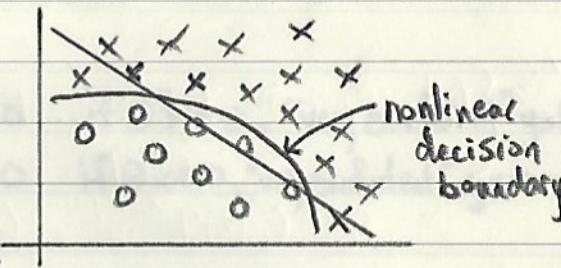
$$b = b - \alpha \frac{\partial \mathcal{L}}{\partial b}$$



## Nonlinear classifiers

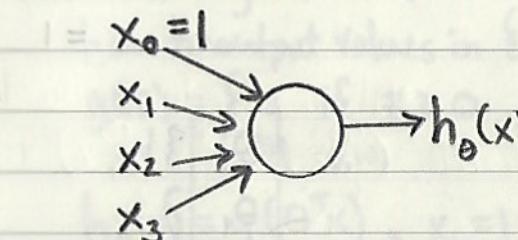
$$\text{logistic regression: } h_\theta(x) = \frac{1}{1+e^{-\theta^T x}}$$

$$\begin{aligned} x | y=1 &\sim \text{Expfamily } (\mu_1) \\ x | y=0 &\sim \text{Expfamily } (\mu_0) \end{aligned} \Rightarrow \text{logistic posterior}$$

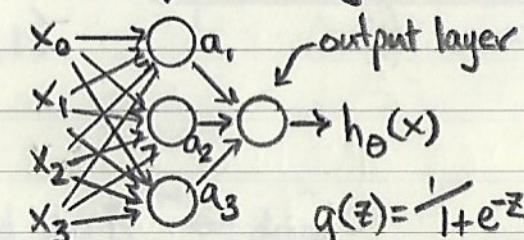


Naive Bayes falls into ExpFamily  $\therefore$  using linear classifier

neural network:



hidden layer



output layer

$$a_1 = g(x^T \theta^{(1)}), a_2 = g(x^T \theta^{(2)}), a_3 = g(x^T \theta^{(3)})$$

$$h_\theta(x) = g(\tilde{a}^T \theta^{(4)}), \tilde{a} = \begin{bmatrix} 1 \\ a_1 \\ a_2 \\ a_3 \end{bmatrix}$$

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n (y^{(i)} - h_\theta(x^{(i)}))^2 \leftarrow \text{quadratic cost function}$$

neural nets not guaranteed to converge to global optimum.

$\uparrow$  non-convex optimization problem

Support Vector Machines (SVM)

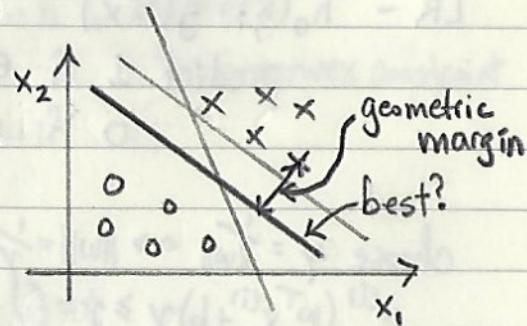
$p(y=1|x; \theta)$  and  $h_\theta(x) = \theta^T x$

compute  $\theta^T x$ , predicts 1 iff  $\theta^T x \geq 0$  if  $\theta^T x > 0$ , "very confident"  $y=1$   
 0 iff  $\theta^T x < 0$  if  $\theta^T x < 0$ , "very confident"  $y=0$

nice if  $\forall i$  s.t.  $y^{(i)}=1$ , have  $\theta^T x \geq 0$

$\forall i$  s.t.  $y^{(i)}=0$ , have  $\theta^T x \leq 0$

assume training set is linearly separable



notation:  $y \in \{-1, 1\}$

have  $h$  output values in  $\{-1, 1\}$

$$g(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ -1 & \text{o.w.} \end{cases}$$

$h_\theta(x) = g(\theta^T x)$ ,  $x_0=1$  and  $x \in \mathbb{R}^{d+1} \leftarrow$  drop

$h_{w,b}(x) = g(w^T x + b)$ ,  $b = \theta_0$ ,  $w = [\theta_1, \dots, \theta_d]^T$ ,  $w \in \mathbb{R}^d$ ,  $x \in \mathbb{R}^d$ ,  $b \in \mathbb{R}$

functional margin of hyperplane  $(w, b)$  w.r.t.  $(x^{(i)}, y^{(i)})$  is

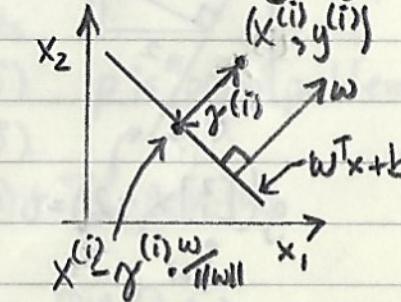
$$\hat{\gamma}^{(i)} = y^{(i)}(w^T x^{(i)} + b) \quad \begin{array}{l} \text{if } y^{(i)}=1 \text{ want } w^T x^{(i)} + b \geq 0 \\ \text{if } y^{(i)}=-1 \text{ want } w^T x^{(i)} + b \leq 0 \end{array}$$

if  $y^{(i)}(w^T x^{(i)} + b) > 0$  then classified  $(x^{(i)}, y^{(i)})$  correctly

$$\hat{\gamma} = \min_i \hat{\gamma}^{(i)}$$
 — functional margin of training set is worst case

normalization condition, like  $\|w\|=1$ , so  $w^T x + b$  is not arbitrarily large

geometric margin —  $\gamma^{(i)}$  (no hat)



$\frac{w}{\|w\|}$  is unit vector normal to  $w^T x + b$

$$w^T(x^{(i)} - \gamma^{(i)} \frac{w}{\|w\|}) + b = 0$$

$$w^T x^{(i)} + b = \gamma^{(i)} \frac{w^T w}{\|w\|} = \gamma^{(i)} \|w\|$$

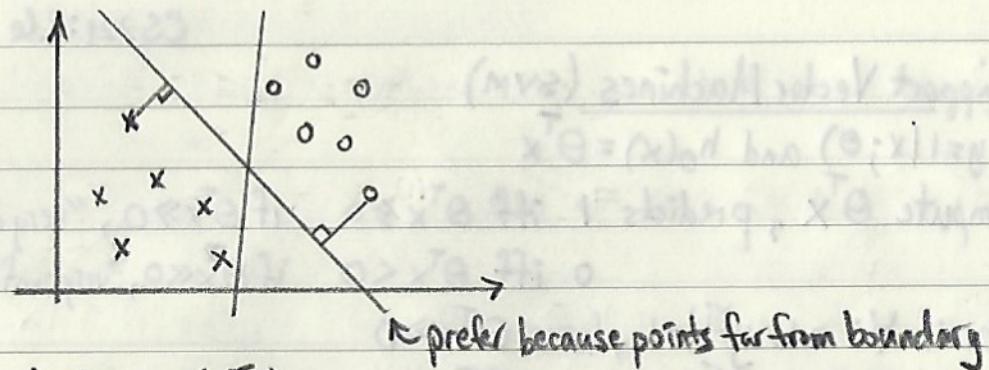
$$\gamma^{(i)} = \left[ \left( \frac{w}{\|w\|} \right)^T x^{(i)} + \frac{b}{\|w\|} \right] y^{(i)}$$

$$\gamma^{(i)} = \frac{\hat{\gamma}^{(i)}}{\|w\|}, \gamma = \min_i \gamma^{(i)}$$

↑ more general to make sign correct

Optimal margin classifier: choose  $w, b$  to maximize  $\gamma$

$$\text{max. margin classifier: } \max_{\gamma, w, b} \gamma \quad \text{s.t. } y^{(i)}(w^T x^{(i)} + b) \geq \gamma \quad \|\omega\|=1$$



$$LR - h_0(x) = g(\theta^T x)$$

predict 1 if  $\theta^T x \geq 0$  want  $\theta^T x \gg 0$   
0 if  $\theta^T x < 0$  want  $\theta^T x \ll 0$

$$\text{choose } \gamma = \frac{1}{\|w\|} \Leftrightarrow \|w\| = \frac{1}{\gamma}$$

$$\Rightarrow y^{(i)}(w^T x^{(i)} + b) \gamma \geq \gamma$$

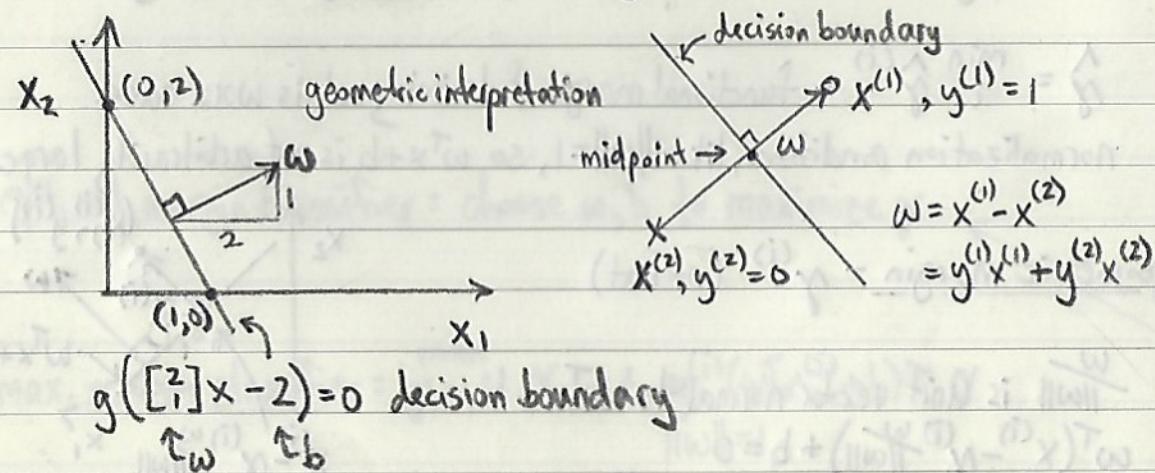
$$\Rightarrow y^{(i)}(w^T x^{(i)} + b) \geq 1, \forall i=1\dots n$$

suppose:  $w = \sum_{i=1}^n \alpha_i x^{(i)}$  - Representer Theorem  
 $= \sum_{i=1}^n \alpha_i y^{(i)} x^{(i)}$   $w$  is a linear combination of examples  
 $\in \{0, 1\}$  and lies in "span" of examples

$$LR - \theta_0 = 0$$

$$SGD - \theta_i = \theta - \alpha (h_0(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$$

$$\text{batch} - \theta_i = \theta - \alpha \sum_{i=1}^n (h_0(x^{(i)}) - y^{(i)}) x^{(i)}$$



## Optimal Margin Classifiers, KKT Conditions, SVM Duals

impose constraint  $\|w\|=1$  or  $|w_i|=1$  or  $w_1^2 + w_2^2 = 1$

↑ can only choose one and then scale to fit  $w^T x + b$

$$\#1 \max_{y, N, b} y \quad \text{s.t. } y^{(i)}(N^T x^{(i)} + b) \geq y, \forall i=1\dots n, \|w\|=1$$

↑ constraint

$$\#2 \max_{\hat{y}, w, b} \hat{y} \quad \text{s.t. } y^{(i)}(w^T x^{(i)} + b) \geq \hat{y}, \forall i=1\dots n$$

↑ non-convex constraint

↑ non-convex objective

impose constraint  $\hat{y}=1$ ,  $\min_i y^{(i)}(w^T x^{(i)} + b) = 1$

$$\#3 \min_{w, b} \frac{1}{2} \|w\|^2 \quad \text{s.t. } y^{(i)}(w^T x^{(i)} + b) \geq 1, i=1\dots n$$

↑ optimal margin classifier

↑ same as  $\max_{w, b} \frac{1}{\|w\|}$



Lagrange multipliers:

$$\min_w f(w) \quad \text{s.t. } h_i(w) = 0, i=1\dots l \rightarrow h(w) = \begin{bmatrix} h_1(w) \\ \vdots \\ h_l(w) \end{bmatrix} = \vec{0}$$

$$\text{Lagrangian: } \mathcal{L}(w, \beta) = f(w) + \sum_i \beta_i h_i(w)$$

↑ Lagrange multipliers

$$\frac{\partial \mathcal{L}}{\partial w} \stackrel{\text{set}}{=} 0, \frac{\partial \mathcal{L}}{\partial \beta} \stackrel{\text{set}}{=} 0$$

for  $w^*$  to be a solution, necessary that  $\exists \beta^* \text{ s.t. } \frac{\partial \mathcal{L}(w^*, \beta^*)}{\partial w} = 0, \frac{\partial \mathcal{L}(w^*, \beta^*)}{\partial \beta} = 0$

↑ there exists

$$\min_w f(w) \quad \text{s.t. } g_i(w) \leq 0, i=1\dots K \quad (g(w) \leq \vec{0}) \quad p \text{ is "primal problem"} \\ h_i(w) = 0, i=1\dots l \quad (h(w) = \vec{0})$$

$$\text{Lagrangian: } \mathcal{L}(w, \alpha, \beta) = f(w) + \sum_{i=1}^K \alpha_i g_i(w) + \sum_{i=1}^l \beta_i h_i(w)$$

$$\text{define: } \Theta_p(w) = \max_{\alpha, \beta} \mathcal{L}(w, \alpha, \beta) \text{ s.t. } \alpha_i \geq 0$$

$$\text{consider: } p^* = \min_w \max_{\alpha, \beta} \mathcal{L}(w, \alpha, \beta) \text{ s.t. } \alpha_i \geq 0 = \min_w \Theta_p(w)$$

↑ value of primal problem

$$\min \frac{1}{2} \|w\|^2 = \frac{1}{2} \left( \sum_{i=1}^n \alpha_i y^{(i)} x^{(i)T} \right) \left( \sum_{j=1}^n \alpha_j y^{(j)} x^{(j)} \right)$$

$$= \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^{(i)} y^{(j)T} \underbrace{x^{(i)T} x^{(j)}}_{\langle x^{(i)}, x^{(j)} \rangle}$$

constraint:

$$y^{(i)} \left( \underbrace{\left( \sum_{j=1}^n \alpha_j y^{(j)} x^{(j)} \right) x^{(i)}}_w + b \right) \geq 1$$

$$y^{(i)} \left( \left( \sum_{j=1}^n \alpha_j y^{(j)} \underbrace{\langle x^{(i)}, x^{(j)} \rangle}_{} \right) + b \right) \geq 1$$

"Kernel trick"

$$\Theta_p(w) \text{ if } g_i(w) > 0$$

$$\text{then } \Theta_p(w) = \infty$$

$$\text{if } h_i(w) \neq 0 \text{ then } \Theta_p(w) = \infty$$

$$\text{otherwise } \Theta_p(w) = f(w)$$

$$\text{thus } \Theta_p(w) = \begin{cases} f(w) & \text{if constraints } g_i, h \text{ are satisfied} \\ +\infty & \text{otherwise} \end{cases}$$

$$\text{so } \min_w \Theta_p(w) = \text{original problem}$$

$$\text{Dual problem} = \Theta_D(\alpha, \beta) = \min_w \mathcal{L}(w, \alpha, \beta)$$

$$d^* = \max_{\alpha \geq 0, \beta} \min_w \mathcal{L}(w, \alpha, \beta) = \max_{\alpha \geq 0, \beta} \Theta_D(\alpha, \beta)$$

$$d^* \leq p^*, \quad \max \min (\dots) \leq \min \max (\dots) \leftarrow \text{true for any function}$$

$$\max_{y \in \{0,1\}} \underbrace{\min_{x \in \{0,1\}} \underbrace{\mathbb{1}\{x=y\}}_{\text{always 0}}} \leq \min_{x \in \{0,1\}} \underbrace{\max_{y \in \{0,1\}} \underbrace{\mathbb{1}\{x=y\}}_{\text{always 1}}}$$

$$\text{Sometimes } d^* = p^*$$

let  $f$  be convex (Hessian is positive semi-definite,  $H \geq 0$ )

suppose  $h_i$  is affine [ $h_i(w) = a_i^T w + b_i$ ] (similar to linear)

and constraint  $g_i$  are (strictly) feasible [ $\exists w \text{ s.t. } \forall i g_i(w) < 0$ ]

then  $\exists w^*, \alpha^*, \beta^*$  s.t.  $w^*$  solves primal

$\alpha^*, \beta^*$  solve dual

$$\text{and } p^* = d^* = \mathcal{L}(w^*, \alpha^*, \beta^*)$$

$$\frac{\partial}{\partial w} \mathcal{L}(w^*, \alpha^*, \beta^*) = 0, \quad \frac{\partial}{\partial \beta} \mathcal{L}(w^*, \alpha^*, \beta^*) = 0$$

$\alpha_i^* g_i(w) = 0 \leftarrow \text{KKT complementarity condition}$

$$g_i(w^*) \leq 0, \quad \alpha_i^* \geq 0$$

Karush-Kuhn-Tucker complementarity condition:

$$\text{if } \alpha_i > 0 \Rightarrow g_i(w^*) = 0$$

$$\text{usually } \alpha_i^* \neq 0 \Leftrightarrow g_i(w^*) = 0$$

$g_i(w)$  is a "active" constraint

Lagrange multipliers:  $\alpha_i, \beta_i \rightarrow \alpha_i$

parameters:  $w \rightarrow w, b$

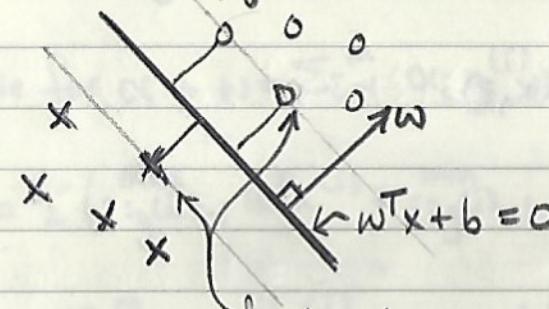
Optimal Margin Classifier:

$$\min \frac{1}{2} \|w\|^2 \text{ s.t. } y^{(i)}(w^T x^{(i)} + b) \geq 1, i = 1 \dots m$$

$$\Leftrightarrow g(w, b) = -y^{(i)}(w^T x^{(i)} + b) + 1 \leq 0$$

$$\alpha_i > 0 \Rightarrow g_i(w, b) = 0 \text{ "active" constraint}$$

$\Leftrightarrow (x^{(i)}, y^{(i)})$  has functional margin = 1



functional margin = 1 (usually  $\alpha_i \neq 0$ )

points called Support Vectors (relatively few)

$\alpha_i = 0$  for non-Support Vectors

$$\mathcal{L}(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^m \alpha_i (y^{(i)}(w^T x^{(i)} + b) - 1) \leftarrow \text{no } \beta_i \text{ because only inequality constraints}$$

Dual problem:  $\Theta_D(\alpha) = \min_{w, b} \mathcal{L}(w, b, \alpha)$

$$\nabla_w \mathcal{L}(w, b, \alpha) = w - \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)} \stackrel{\text{set}}{=} 0 \Rightarrow w = \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)}$$

$$\frac{\partial \mathcal{L}}{\partial b} = - \sum_{i=1}^m \alpha_i y^{(i)} \stackrel{\text{set}}{=} 0$$

$$\max \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_i \sum_j y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle \text{ s.t. } \alpha_i \geq 0$$

$$\sum y^{(i)} \alpha_i = 0$$

solve for  $\alpha_i, b$

$$\begin{aligned} \text{prediction: } h_{w,b}(x) &= g(w^T x + b) \\ &= g((\sum_i \alpha_i y^{(i)} x^{(i)})^T x + b) \\ &= g((\sum_i \alpha_i y^{(i)} \langle x^{(i)}, x \rangle) + b) \end{aligned}$$

### Kernel trick

1) write algo in terms of  $\langle x^{(i)}, x^{(j)} \rangle$

$$\uparrow_x \uparrow_z \langle x, z \rangle$$

2) mapping  $x \mapsto \phi(x)$

$\uparrow$  very high dimensional

3) compute  $K(x, z) = \phi(x)^T \phi(z)$

4) replace  $\langle x, z \rangle$  in algo by  $K(x, z)$

$$\mathcal{L} = \frac{1}{2} w^T w - \sum_{i=1}^m \alpha_i (y^{(i)} (w^T x^{(i)}) + b) - 1$$

$$w^T w = (\sum_{i=1}^m \alpha_i y^{(i)} x^{(i)})^T (\sum_{j=1}^m \alpha_j y^{(j)} x^{(j)}) \leftarrow \text{from previous}$$

$$\begin{aligned} \mathcal{L} &= \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle - \sum_{i=1}^m \sum_{j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle \\ &\quad + \sum_{i=1}^m \alpha_i \end{aligned}$$

$$= \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle = W(\alpha)$$

$$\text{Dual problem: } \max W(\alpha) \text{ s.t. } \alpha_i \geq 0, \sum_i y_i \alpha_i = 0$$

if  $\sum_i y_i \alpha_i \neq 0$  then  $\Theta_D(\alpha) = -\infty$

if  $\sum_i y_i \alpha_i = 0$  then  $\Theta_D(\alpha) = W(\alpha)$

$$\text{Solve for } \alpha \rightarrow w = \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)}$$

$$b = \frac{1}{2} \left( \max_{i: y^{(i)}=-1} w^T x^{(i)} + \min_{i: y^{(i)}=1} w^T x^{(i)} \right)$$

$$w = \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)}$$

$$h_{w,b}(x) = g(w^T x + b)$$

$$w^T x + b = \sum_{i=1}^m \alpha_i y^{(i)} \langle x^{(i)}, x \rangle + b \leftarrow \text{only need to compute for support vectors } (\alpha_i \neq 0)$$

Kernels -  $x^{(i)}$  very high dimensional ( $x^{(i)} \in \mathbb{R}^\infty$ )

$\langle x^{(i)}, x^{(j)} \rangle$  efficiently computed

convex optimization problem :  $\min \frac{1}{2} \|w\|^2$   
 s.t.  $y^{(i)}(w^T x^{(i)} + b) \geq 1, i = 1 \dots m$   
 assumes data linearly separable

dual problem :  $\max \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j y^{(i)} y^{(j)} \alpha_i \alpha_j \langle \underbrace{x^{(i)}, x^{(j)}}_{(x^{(i)})^T x^{(j)}} \rangle$   
 s.t.  $\alpha_i \geq 0, \sum_i y_i \alpha_i = 0$

$w = \sum_i \alpha_i y^{(i)} x^{(i)}$ ,  $h_{w,b}(x) = g(w^T x + b)$  ← make prediction  
 $= g(\sum_i \alpha_i y^{(i)} \langle \underbrace{x^{(i)}, x}_{(x^{(i)})^T x} \rangle + b)$

have  $x \in \mathbb{R}$  - living area of house

$x \xrightarrow{\phi} \begin{bmatrix} x \\ x^2 \\ x^3 \\ x^4 \end{bmatrix} \rightarrow \phi(x)$ , replace  $\langle x^{(i)}, x^{(j)} \rangle$  with  $\langle \phi(x^{(i)}), \phi(x^{(j)}) \rangle$   
 ↑ high dimensional, possibly infinite

Kernel function :  $K(x^{(i)}, x^{(j)}) = \langle \phi(x^{(i)}, x^{(j)}) \rangle$  ← replace all inner products

$$\begin{aligned}
 x, z \in \mathbb{R}^n \\
 K(x, z) &= (x^T z)^2 = (\sum_{i=1}^d x_i z_i)(\sum_{j=1}^d x_j z_j) = \sum_{i=1}^d \sum_{j=1}^d (x_i x_j)(z_i z_j) \\
 &= (\phi(x))^T (\phi(z))
 \end{aligned}$$

$\phi(x) = \begin{bmatrix} x_1 x_1 \\ x_1 x_2 \\ x_1 x_3 \\ x_2 x_1 \\ x_2 x_2 \\ x_2 x_3 \\ x_3 x_1 \\ x_3 x_2 \\ x_3 x_3 \end{bmatrix}$   
 $O(d^2)$  to compute  $\phi(x)$   
 $O(d)$  to compute  $K(x, z)$

$$\begin{aligned}
 n=3 \rightarrow & \begin{bmatrix} x_1 x_1 \\ x_1 x_2 \\ x_1 x_3 \\ x_2 x_1 \\ x_2 x_2 \\ x_2 x_3 \\ x_3 x_1 \\ x_3 x_2 \\ x_3 x_3 \end{bmatrix} \downarrow^2 \\
 & \begin{bmatrix} \sqrt{2}c x_1 \\ \sqrt{2}c x_2 \\ \sqrt{2}c x_3 \\ c \end{bmatrix} \leftarrow K(x, z) = (x^T z + c)^2
 \end{aligned}$$

$$x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \in \mathbb{R}^d \quad \phi(x) = \begin{bmatrix} x_1 x_1 \\ x_1 x_2 \\ x_1 x_3 \\ x_2 x_1 \\ x_2 x_2 \\ x_2 x_3 \\ x_3 x_1 \\ x_3 x_2 \\ x_3 x_3 \end{bmatrix} \in \mathbb{R}^{d^2} \quad \phi(z) = \begin{bmatrix} z_1 z_1 \\ z_1 z_2 \\ z_1 z_3 \\ z_2 z_1 \\ z_2 z_2 \\ z_2 z_3 \\ z_3 z_1 \\ z_3 z_2 \\ z_3 z_3 \end{bmatrix} \in \mathbb{R}^{d^2}$$

if  $x, z$  are similar,  $K(x, z) = \phi(x)^T \phi(z)$  is "large"  
 " dissimilar " " " " "small"

$$K(x, z) = \exp\left(-\frac{(x-z)^2}{2\sigma^2}\right)$$

$$K(x, x) = \phi(x)^T \phi(x) \geq 0$$

more generally  $x^{(1)}, \dots, x^{(t)}$

$$\begin{aligned} & (\sum_{i=1}^t z_i x^{(i)})^T (\sum_{j=1}^t z_j x^{(j)}) \geq 0 \\ & = \sum_{i=1}^t \sum_{j=1}^t z_i z_j K(x^{(i)}, x^{(j)}) \geq 0 \\ & = z^T K z \geq 0 \quad \forall z \in \mathbb{R}^t \end{aligned}$$

$K(x, z) = (x^T z + c)^t \rightarrow \binom{n+t}{t}$  features of all monomials up to degree  $t$   
 $\uparrow n+t$  choose  $t \approx (d+t)^t$

$$x \mapsto \phi(x), z \mapsto \phi(z)$$

$\uparrow$  attributes  $\uparrow$  feature vector

$$\langle \phi(x), \phi(z) \rangle$$

$K(x, z)$  - large if  $x, z$  similar

small if  $x, z$  dissimilar

$$K(x, z) = \exp\left(-\frac{\|x-z\|^2}{2\sigma^2}\right) - \text{Gaussian Kernel}$$

is  $K$  a valid Kernel?  $\exists \phi$  s.t.  $K(x, z) = \langle \phi(x), \phi(z) \rangle$ ?

Suppose  $K$  is a Kernel, let  $\{x^{(1)}, \dots, x^{(n)}\}$  be given

$$\text{let } K \in \mathbb{R}^{n \times n}, K_{ij} = K(x^{(i)}, x^{(j)})$$

$\uparrow$  kernel matrix

$$\begin{aligned} \text{then for any vector } z \in \mathbb{R}^n, z^T K z &= \sum_i \sum_j z_i K_{ij} z_j \\ &= \sum_i \sum_j z_i \phi(x^{(i)})^T \phi(x^{(j)}) z_j \\ &= \sum_i \sum_j z_i \sum_k (\phi(x^{(i)})_k) (\phi(x^{(j)})_k) z_j \\ &= \sum_k \sum_i \sum_j z_i (\phi(x^{(i)})_k) (\phi(x^{(j)})_k) z_j \\ &= \sum_k (\sum_i z_i \phi(x^{(i)})_k)^2 \geq 0 \end{aligned}$$

Kernel positive semidefinite  $\uparrow$

valid  $\downarrow$

Theorem (Mercer): Let  $K(x, z)$  be given, then  $K$  is a valid (Mercer) Kernel

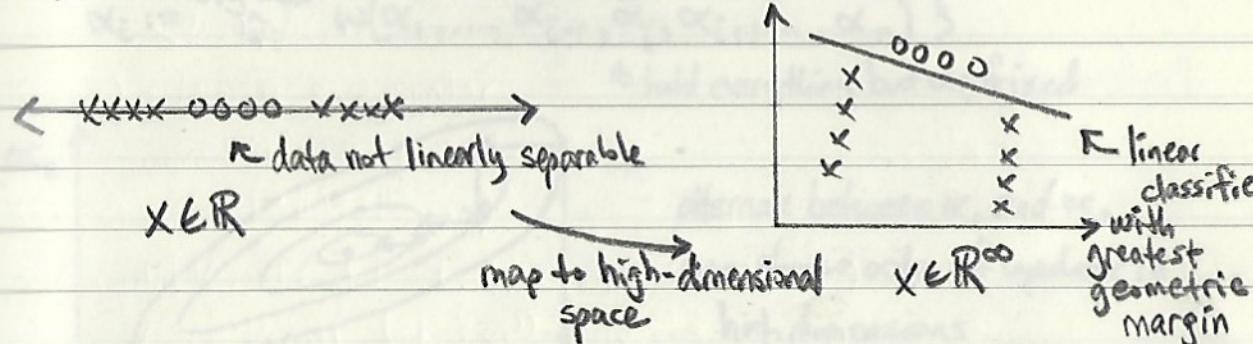
(i.e.  $\exists \phi$  s.t.  $K(x, z) = \phi(x)^T \phi(z)$ ) if and only if

for all  $\{x^{(1)}, \dots, x^{(t)}\}$  ( $t < \infty$ ) the Kernel matrix  $K \in \mathbb{R}^{t \times t}$  is symmetric positive semidefinite.

$K(x, x) = 1 \neq \phi(x)^T \phi(x)$  - inner product of vector with itself is  $\geq 0$

choose  $K(x, z) = \exp\left(-\frac{\|x-z\|^2}{2\sigma^2}\right)$  Gaussian -  $\infty$  dimensional feature space  
 or  $(x^T z + c)^d$

replace  $\langle x^{(i)}, x^{(j)} \rangle$  with  $K(x^{(i)}, x^{(j)})$ ,  $x^{(i)} \rightarrow \phi(x^{(i)})$



Kernels much broader than SVMs - most linear algorithms can be rewritten using  $\langle x^{(i)}, x^{(j)} \rangle \rightsquigarrow K(x^{(i)}, x^{(j)})$

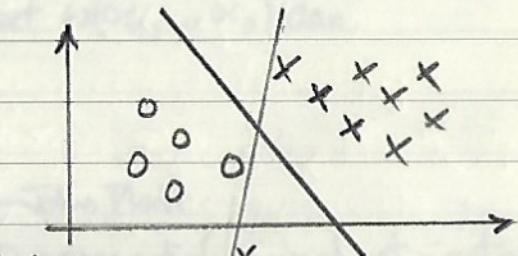
L<sub>1</sub> norm soft margin SVM:

$$\min \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i$$

penalty term

$$\text{s.t. } y^{(i)}(w^T x^{(i)} + b) \geq 1 - \xi_i, \xi_i \geq 0 \forall i=1, \dots, n$$

functional margin



one outlier skewing decision boundary

if  $y^{(i)}(w^T x^{(i)} + b) > 0 \Rightarrow$  classified correctly

$$\mathcal{L}(w, b, \xi, \alpha, r) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i - \sum_{i=1}^n \alpha_i (y^{(i)}(w^T x^{(i)} + b) - 1 + \xi_i) - \sum_{i=1}^n \alpha_i \xi_i$$

Lagrangian

hyperparameter

$$\max_w w(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_i \sum_j y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle$$

s.t.  $\sum_{i=1}^n y^{(i)} \alpha_i = 0, 0 \leq \alpha_i \leq C \forall i=1, \dots, n$

$$\alpha_i = 0 \Rightarrow y^{(i)}(w^T x^{(i)} + b) \geq 1$$

$$\alpha_i = C \Rightarrow y^{(i)}(w^T x^{(i)} + b) \leq 1$$

$$0 < \alpha_i < C \Rightarrow y^{(i)}(w^T x^{(i)} + b) = 1$$

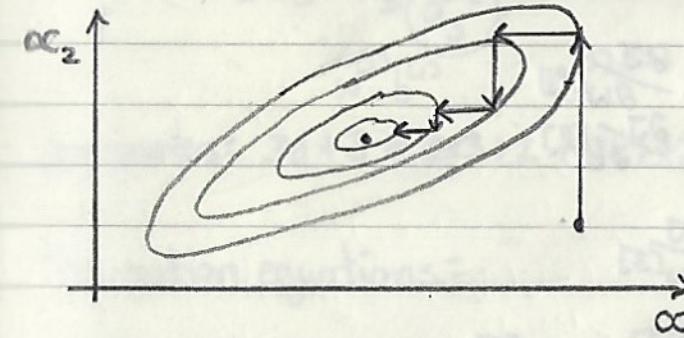
### Coordinate Ascent

$\max(\alpha_1, \dots, \alpha_n)$  - no constraints on  $\alpha_i$ 's

Repeat { for  $i = 1$  to  $n$

$$\alpha_i := \arg \max_{\alpha_i} w(\alpha_1, \dots, \alpha_{i-1}, \hat{\alpha}_i, \alpha_{i+1}, \dots, \alpha_n)$$

↑ hold everything but  $\alpha_i$  fixed



alternate between  $\alpha_1$  and  $\alpha_2$ ,  
can choose order of update in  
high dimensions

takes more steps than Newton's method, but  $w(\alpha_1, \dots, \alpha_n)$  can be inexpensive

$$\sum_{i=1}^n y^{(i)} \alpha_i = 0 \text{ - constraint}$$

chang 2  $\alpha_i$ 's at a time

select  $\alpha_i, \alpha_j$

hold all  $\alpha_k$  fixed except  $\alpha_i, \alpha_j$

optimize  $w(\alpha)$  w.r.t.  $\alpha_i, \alpha_j$  s.t. constraints

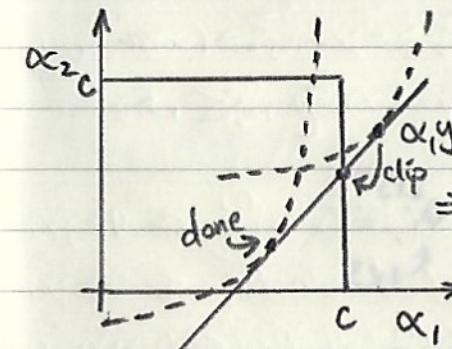
John Platt

SMO - sequential minimal optimization

update =  $\alpha_1$  and  $\alpha_2$  (select 2  $\alpha_i$  on every iteration)

$$\alpha_1 y^{(1)} + \alpha_2 y^{(2)} = - \sum_{i=2}^n \alpha_i y^{(i)} = \xi \text{ - constraint}$$

$0 \leq \alpha_i \leq C$  - box constraint



$$w(\alpha_1, \dots, \alpha_n) =$$

$$w(y^{(1)}(\xi - \alpha_2 y^{(2)}, \alpha_2, \dots) =$$

$$\Rightarrow \alpha_1 = \frac{1}{y^{(1)}} (\xi - \alpha_2 y^{(2)})$$

$$a\alpha_2^2 + b\alpha_2 + c$$

Loss function:

$$J(\hat{y}, y) = \frac{1}{n} \sum_{i=1}^n L^{(i)}$$

$$L^{(i)} = -[y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log (1-\hat{y}^{(i)})]$$

Backward propagation:

$$\forall l=1\dots 3 \quad w^{[l]} := w^{[l]} - \alpha \frac{\partial J}{\partial w^{[l]}} \\ b^{[l]} := b^{[l]} - \alpha \frac{\partial J}{\partial b^{[l]}}$$

$$\frac{\partial J}{\partial w^{[3]}} = \underbrace{\frac{\partial J}{\partial a^{[3]}}}_{\substack{\text{from loss} \\ \text{function}}} \cdot \underbrace{\frac{\partial a^{[3]}}{\partial z^{[3]}}}_{\substack{\text{activation} \\ \text{function}}} \cdot \underbrace{\frac{\partial z^{[3]}}{\partial w^{[3]}}}_{\substack{\text{gradient} \\ \text{of weights}}}$$

$$\frac{\partial J}{\partial w^{[2]}} = \underbrace{\frac{\partial J}{\partial z^{[3]}}}_{\substack{\text{from loss} \\ \text{function}}} \cdot \underbrace{\frac{\partial z^{[3]}}{\partial a^{[2]}}}_{\substack{\text{activation} \\ \text{function}}} \cdot \underbrace{\frac{\partial a^{[2]}}{\partial z^{[2]}}}_{\substack{\text{gradient} \\ \text{of weights}}} \cdot \underbrace{\frac{\partial z^{[2]}}{\partial w^{[2]}}}_{\substack{\text{gradient} \\ \text{of weights}}}$$

$$\frac{\partial J}{\partial w^{[1]}} = \underbrace{\frac{\partial J}{\partial z^{[2]}}}_{\substack{\text{from loss} \\ \text{function}}} \cdot \underbrace{\frac{\partial z^{[2]}}{\partial a^{[1]}}}_{\substack{\text{activation} \\ \text{function}}} \cdot \underbrace{\frac{\partial a^{[1]}}{\partial z^{[1]}}}_{\substack{\text{gradient} \\ \text{of weights}}} \cdot \underbrace{\frac{\partial z^{[1]}}{\partial w^{[1]}}}_{\substack{\text{gradient} \\ \text{of weights}}}$$

$$\begin{aligned} \frac{\partial L^{(i)}}{\partial w^{[3]}} &= - \left[ y^{(i)} \underbrace{\frac{\partial}{\partial w^{[3]}} \log \sigma(w^{[3]} a^{[2]} + b^{[3]})}_{\substack{\text{activation} \\ \text{function}}} + (1-y^{(i)}) \underbrace{\frac{\partial}{\partial w^{[3]}} \log (1-\sigma(w^{[3]} a^{[2]} + b^{[3]}))}_{\substack{\text{activation} \\ \text{function}}} \right] \\ &= - \left[ y^{(i)} \cdot \frac{1}{\sigma(a^{[2]})} (1-\sigma(a^{[2]})) a^{[2]T} + (1-y^{(i)}) \frac{1}{1-\sigma(a^{[2]})} (-1) \sigma(a^{[2]}) (1-\sigma(a^{[2]})) a^{[2]T} \right] \\ &= - \left[ y^{(i)} (1-\sigma(a^{[2]})) a^{[2]T} - (1-y^{(i)}) \sigma(a^{[2]}) a^{[2]T} \right] \\ &= - \left[ y^{(i)} a^{[2]T} - \sigma(a^{[2]}) a^{[2]T} \right] \\ &= - (y^{(i)} - \sigma(a^{[2]})) a^{[2]T} \end{aligned}$$

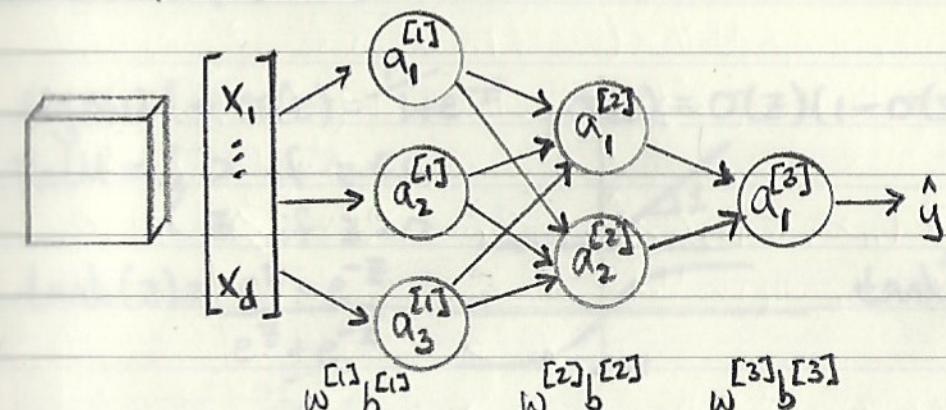
$$\frac{\partial L^{(i)}}{\partial w^{[2]}} = \underbrace{\frac{\partial}{\partial a^{[3]}}}_{\substack{\text{gradient} \\ \text{of activation} \\ \text{function}}} \cdot \underbrace{\frac{\partial a^{[3]}}{\partial z^{[3]}}}_{\substack{\text{activation} \\ \text{function}}} \cdot \underbrace{\frac{\partial z^{[3]}}{\partial a^{[2]}}}_{\substack{\text{gradient} \\ \text{of weights}}} \cdot \underbrace{\frac{\partial a^{[2]}}{\partial z^{[2]}}}_{\substack{\text{gradient} \\ \text{of activation} \\ \text{function}}} \cdot \underbrace{\frac{\partial z^{[2]}}{\partial w^{[2]}}}_{\substack{\text{gradient} \\ \text{of weights}}}$$

$$\frac{\partial L^{(i)}}{\partial w^{[3]}} = \underbrace{\frac{\partial}{\partial z^{[3]}}}_{\substack{\text{gradient} \\ \text{of activation} \\ \text{function}}} \cdot \underbrace{\frac{\partial z^{[3]}}{\partial w^{[3]}}}_{\substack{\text{gradient} \\ \text{of weights}}} \\ - (y^{(i)} - \sigma(a^{[2]})) \uparrow \quad \uparrow a^{[2]T}$$

$$\frac{\partial L^{(i)}}{\partial w^{[2]}} = w^{[3]T} * \underbrace{a^{[2]} (1-\sigma(a^{[2]}))}_{\substack{\text{gradient} \\ \text{of activation} \\ \text{function}}} (\sigma(a^{[2]}) - y^{(i)}) a^{[1]T}$$

## Machine Learning Neural Networks

CS229 : L9



$$\text{parameters: } 3d + 3 + 2 \times 3 + 2 + 2 \times 1 + 1$$

Propagation equations:

$$z^{[1]} = w^{[1]} x + b^{[1]}$$

$$z^{[2]} = w^{[2]} a^{[1]} + b^{[2]}$$

$$z^{[3]} = w^{[3]} a^{[2]} + b^{[3]}$$

$$a^{[1]} = \sigma(z^{[1]})$$

$$a^{[2]} = \sigma(z^{[2]})$$

$$a^{[3]} = \sigma(z^{[3]})$$

$$z^{[1]} = \underbrace{\begin{bmatrix} 1 & \dots & 1 \end{bmatrix}}_{\substack{\text{examples} \\ d}} \uparrow$$

$$x = \begin{bmatrix} x^{(1)} & \dots & x^{(n)} \end{bmatrix} \uparrow$$

$$z^{[1]} = \begin{bmatrix} 1 & \dots & 1 \end{bmatrix} \uparrow$$

$$z^{[1]} = w^{[1]} x + b^{[1]}$$

$$z^{[2]} = w^{[2]} a^{[1]} + b^{[2]}$$

$$z^{[3]} = w^{[3]} a^{[2]} + b^{[3]}$$

$$a^{[1]} = \sigma(z^{[1]})$$

$$a^{[2]} = \sigma(z^{[2]})$$

$$a^{[3]} = \sigma(z^{[3]})$$

$$z^{[1]} = \underbrace{\begin{bmatrix} 1 & \dots & 1 \end{bmatrix}}_{\substack{\text{examples} \\ d}} \uparrow$$

$$z^{[2]} = w^{[2]} a^{[1]} + b^{[2]}$$

$$z^{[3]} = w^{[3]} a^{[2]} + b^{[3]}$$

$$a^{[1]} = \sigma(z^{[1]})$$

$$a^{[2]} = \sigma(z^{[2]})$$

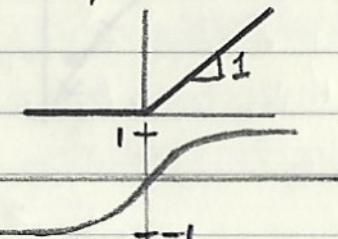
$$a^{[3]} = \sigma(z^{[3]})$$

## Activation functions:

$$\text{Sigmoid} = \sigma(z) = \frac{1}{1+e^{-z}}, \sigma'(z) = \sigma(z)(1-\sigma(z))$$

$$\text{ReLU} = \begin{cases} 0 & \text{if } z \leq 0 \\ z & \text{if } z > 0 \end{cases}$$

$$\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$



$$\tanh'(z) = (-\tanh(z))^2$$

Sigmoid and tanh suffer from vanishing gradient problem (not ReLU)

with no activation neural network reduces to linear regression

typically use same activation for each layer

Initialization:

$$\begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix} = \frac{1}{n} \sum_{i=1}^n x^{(i)} \quad \leftarrow \text{center data around origin}$$

$$x = x - \mu$$

$$\begin{aligned} \Sigma &= \frac{1}{n} \sum_{i=1}^n x^{(i)T} x^{(i)} && \leftarrow \text{normalize} \\ x &= \Sigma^{-\frac{1}{2}} x \end{aligned}$$

compute  $\Sigma$  and  $\mu$  once for training data and reuse for test data

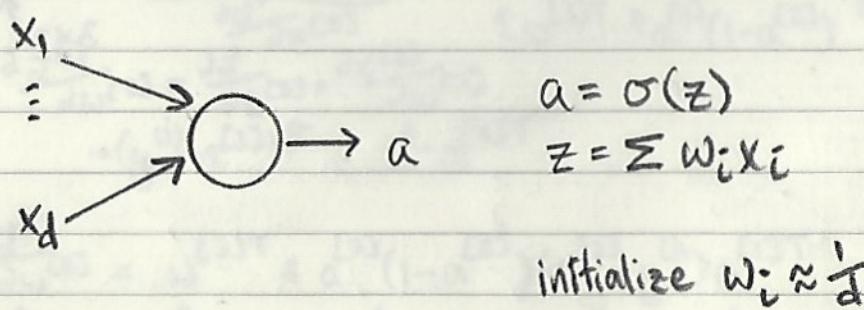
exploding gradient:

$$\hat{y} = w^{[L]} w^{[L-1]} \dots w^{[1]} \Rightarrow w^{[L]} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}, \hat{y} = \begin{bmatrix} 1.5^L & 0 \\ 0 & 1.5^L \end{bmatrix}$$

$$\hat{y} = \begin{bmatrix} 0.5^L & 0 \\ 0 & 0.5^L \end{bmatrix}$$

$$w \approx \sqrt{\frac{1}{d^{[L]}}} \quad \text{- works well for sigmoid}$$

$$w \approx \sqrt{\frac{2}{d^{[L]}}} \quad \text{- " ReLU (He initialization)}$$

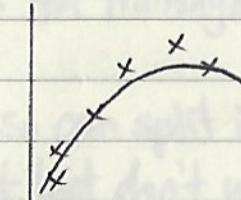


## Machine Learning

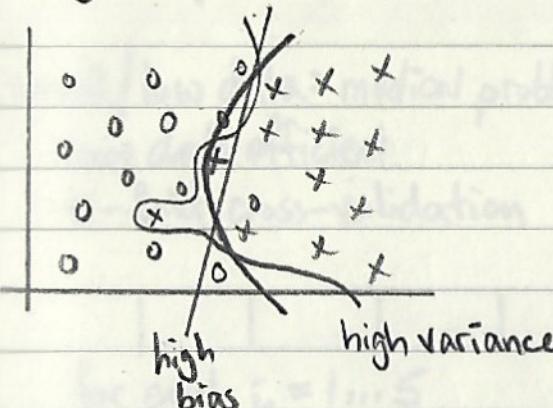
CS229-L10

Bias-Variance, Cross Validation, Model Selection, Regularization

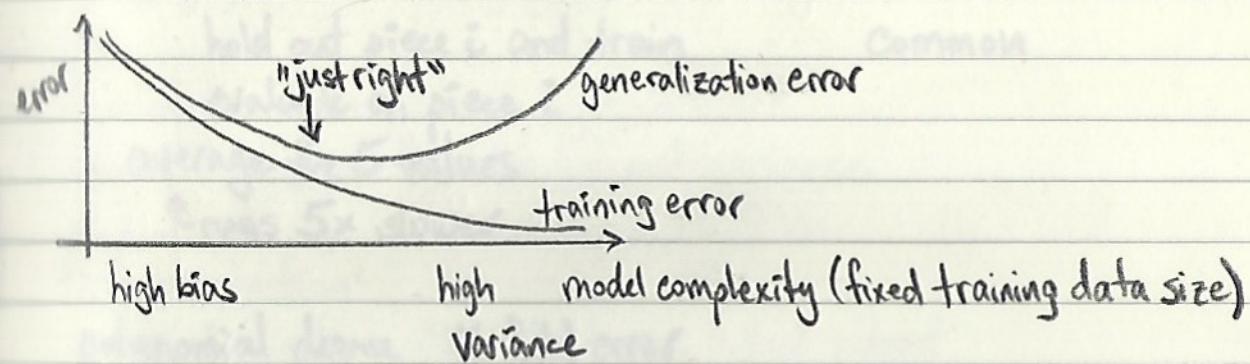
high bias



high variance



high bias



model selection:

$$\Theta_0 + \Theta_1 x$$

linear

$$\Theta_0 + \Theta_1 x + \Theta_2 x^2$$

quadratic

$$\Theta_0 + \dots + \Theta_3 x^3$$

cubic

Hold-out cross validation: split data

70%	30%
training	hold-out CV

- train on training data
- evaluate on dev set
- pick best model on dev set
- retrain on 100% of data (optional)

With very large data sets, dev % is much smaller  
 ↳ can be 99% / 1% for training/dev

Can overfit to dev set, so can split into train/dev/test  
 evaluate with test set but don't use to make any decisions  
 Often skip test set in industry - not publishing results

small / low data: medical problems  
 more data efficient  
 K-fold cross-validation



for each  $i = 1 \dots 5$

hold out piece  $i$  and train  
 evaluate on piece  $i$   
 average the 5 values

↳ runs 5x slower

5-fold CV

↳ 10-fold CV is most

common

polynomial degree      K-fold error

1	3.5
2	2.7
3	3.8
4	5.0

leave-one-out CV =  $K = n$

use for tiny ( $< 100$ ) examples

can use K-fold on training while reserving test for evaluation

target  $\sim 1\%$  error

train 10% error } high bias - does poorly on train  
CV 10.5% error } and CV

train 1% } high variance - overfitting training data  
CV 10%

train 0.5% } success!  
CV 0.6%

train 10% } high bias and variance - can happen in very  
CV 22% } high-dimensional spaces

more data only helps with high variance  
more features can fix bias but not variance

### Feature selection:

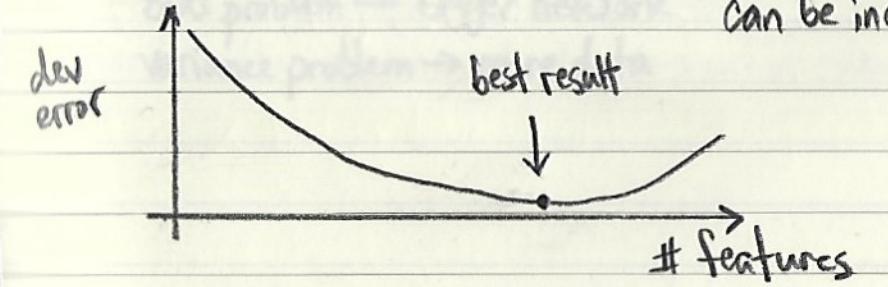
very high-dimensional features - select subset  
forward search: start with  $\mathcal{F} = \{\}$  (no features)

repeat & try adding each feature to  $\mathcal{F}$   
train model for each addition

add best choice to  $\mathcal{F}$  as measured by CV }

difficult to find features that interact

choose set of features that works best on dev set  
can be inefficient



Backward search: start with all features and subtract  
 "Wrapper" feature selection  
 ↗ learning algo called as subroutine

"Filter" feature selection:

for each feature  $i$  compute a score that tells us how informative  $x_i$  is about  $y$   
 pick top  $K$  features

sort features based on  $|\text{corr}(x_i, y)|$

$x_i$	0	1
word	0.1	0.3
label	0.1	0.5

$$\text{MI}(x_i, y) = \sum_{x_i \in \{0, 1\}} \sum_{y \in \{0, 1\}} p(x_i, y) \log \frac{p(x_i, y)}{p(x_i)p(y)}$$

$$= \text{difference}(p(x_i, y), p(x_i)p(y))$$

↑ KL divergence

pick top  $K$  features - pick " $K$ " based on dev set performance

Regularization: reduces variance

$$\min_{\Theta} \frac{1}{2} \sum_{i=1}^n \|y^{(i)} - \Theta^T x^{(i)}\|^2 + \lambda \|\Theta\|^2 \quad \text{linear regression}$$

$$\max_{\Theta} \sum_{i=1}^n \log p(y^{(i)} | x^{(i)}; \Theta) - \underbrace{\lambda \|\Theta\|^2}_{\text{regularization}} \quad \text{logistic regression}$$

penalize large values of  $\Theta$  - smooths out function

deep neural networks diminish bias-variance trade off  
 bias problem → bigger network  
 variance problem → more data

# Machine Learning

## Practical Advice for ML Projects

7 steps overview:

- 1) Acquire data
- 2) look at data - after every step
- 3) Create train/dev/test splits
- 4) Create refine specification
- 5) build model - simplest that works
- 6) measurement
- 7) repeat

Cold-start = feature doesn't exist yet } already incorporating  
legal/ethical issues to look at data } bias

data artifacts are hard - models may pick up unintended signal

"Become one with the data." - Karpathy  
have the right people look at the data

train/dev/test - randomly sampled

test is proxy for real world

leakage - (nearly) same example in train and dev

  ^ performance overstated

dev for tuning parameters

a good specification has little ambiguity

specification must be embodied is a set of examples - test set!

inner annotator agreement - how often do they agree?  
meaningless accuracy problem  
train annotators

subtle problem: consistency in test sets  
" : spec creep - needs to be consistent

avoid getting bogged down in models  
build simpler models even after fancy models  
how accurate is the baseline

Ablation studies - remove one feature at a time from top model

measure end-to-end quality metrics  
measure simple things

Slice-based monitoring - overall performance might be less important than given "slice"  
record and scoreboard on slices  
monitoring should support fine-grained reporting

Avoid unknown mistakes - popularity shift / cold start  
- input shift

retrain models periodically  
labels and input drift over time

a well-running ML system is a rewritten poorly running system

errors in specification =

class schema issue - two distinguishable classes merged  
- one class split into two

unknowable class - information not available to model

unrealized structure

test set label variance

change between test set versions

class confusion matrix

"ground truth" is constructed - fix the specification  
fix the data

measure error

select more labels - uniform random sampling  
importance-based sampling

error bucketing is still critical

add labels to drive model to predict the right class

model diagnostics = try simple models first  
linear or logistic regression with simple features

H.1: PS523

metre dinner along collision or 2nd order diff gradient-like

= gradient descent in 2nd order

beginning of learning rate = cut - mean error's rank  
cut size flip each time -

failure of additive bias minimization - 2nd order add-on terms

multiple local optima

minimize total loss

minimize the best model's graph

various nice things

gradient descent with step - bottleneck in "short training"

step with step - gradient descent

done gradient

gradient descent minima - global from tools  
gradients, local - nothinglocal like in gradient descent  
local like with taking of both min at global like

not global despite gradient = gradient paths taken

2nd order stepsize like minimum staged to until

train vs. test error:

if error too high, model needs more capacity

model needs more data or less complexity

train and test should be close

variance diagnostic: sample data set

k-fold cross validation

calibration plots: bump means there's a lurking class

$$\mu^{(t)} = \arg \min Z(\mu^{(t)} - X)$$

repeat until no changes

randomly select  $\mu^{(t)}$  from  $\mu^{(0)}$

for  $i = 1$  to  $n$ :

$C_i^{(t)} = \text{min}_{\mu^{(t)}} \|x_i - \mu^{(t)}\|$

for each  $i = 1$  to  $n$ :

$\mu^{(t+1)} = \text{min}_{\mu^{(t)}} \sum_{i=1}^n C_i^{(t)} \delta_i(\mu^{(t)})$  ← compute centroid

$$J(C, \mu) = \sum_{i=1}^n \|x_i - \mu\|_2^2$$

find and find minimum - it gets stuck on wrong side of cluster

minimization is NP-hard

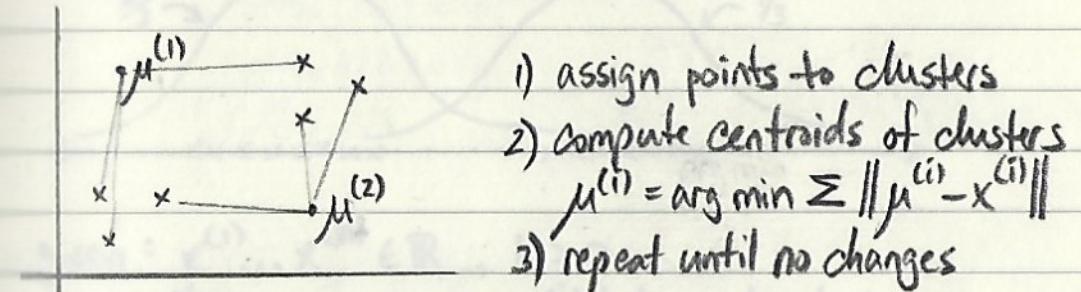
K-means++ (2007) -  $\mu$  initializationand broadcast data to pick  $K$

## Machine Learning

K-means, GMM, Expectation Maximization

unsupervised learning - no labels  
- clustering

given =  $x^{(1)} \dots x^{(n)}$ ,  $x^{(i)} \in \mathbb{R}^d$ ,  $K = \# \text{ of clusters}$   
find assignment of  $x^{(i)} \rightarrow \text{clusters}$



randomly select  $\mu^{(1)} \dots \mu^{(K)} \in \mathbb{R}^d$

for  $i = 1 \dots n$ :

$$c^{(i)} = \arg \min_{j \in 1 \dots k} \| x^{(i)} - \mu_j \|^2$$

for each  $j = 1 \dots k$ :

$$\mu^{(j)} = \frac{1}{|S_j|} \sum_{i \in S_j} x^{(i)}, S_j = \{i : c^{(i)} = j\}$$

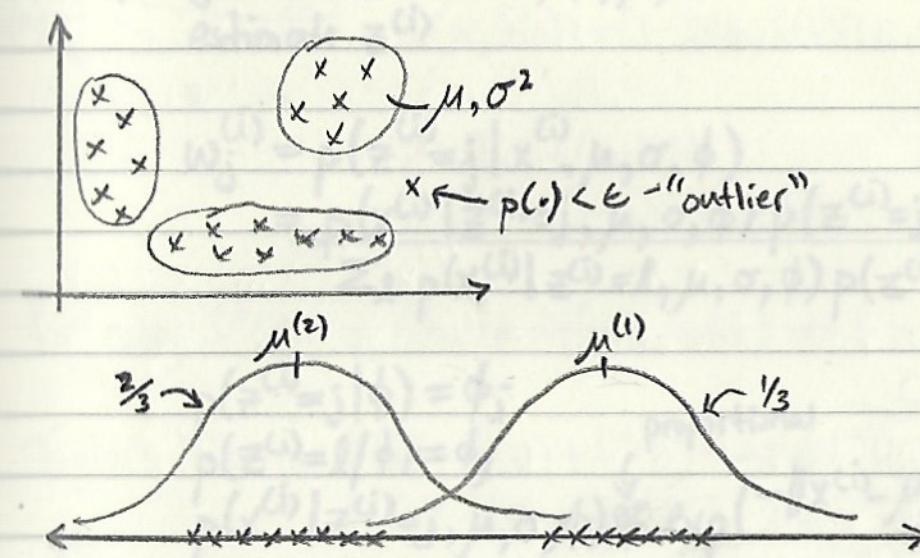
$$J(c, \mu) = \sum_{i=1}^n \| x^{(i)} - \mu_{c^{(i)}} \|^2$$

may not find minimum -  $\mu$  gets stuck on wrong side of cluster  
minimization is NP-hard

K-means++ (2007) -  $\mu$  initialization

need to understand data to pick  $k$

## Mixture of Gaussians:



given:  $x^{(1)}, \dots, x^{(n)} \in \mathbb{R}$ ,  $K > 0$

$p(z^{(i)}=j)$  - probability  $x^{(i)}$  belongs to cluster  $j$

## Gaussian Mixture Model (GMM):

$$p(x^{(i)}, z^{(i)}) = p(x^{(i)} | z^{(i)}) p(z^{(i)})$$

$$z^{(i)} \sim \text{multinomial}(\phi) \quad \sum_{j=1}^K \phi_j = 1, \phi_j \geq 0$$

$$x^{(i)} | z^{(i)}=j \sim N(\mu_j, \sigma_j^2)$$

estimate  $\phi, \mu_j, \sigma_j^2$ ,  $z^{(i)}$  is hidden/latent - postulate

(e-step) 1. guess values of  $z^{(i)}$  for each point

(m-step) 2. update  $\mu_j, \sigma_j, \phi$  parameters

e-step: given  $x^{(1)}, \dots, x^{(n)}$ ,  $\phi, \mu, \sigma$   
estimate  $z^{(i)}$

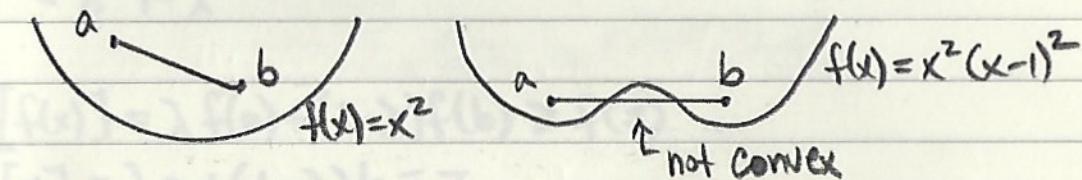
$$\begin{aligned} w_j^{(i)} &= p(z^{(i)}=j | x^{(i)}, \mu, \sigma, \phi) \\ &= \frac{p(x^{(i)} | z^{(i)}=j, \mu, \sigma, \phi) p(z^{(i)}=j | \phi)}{\sum_l p(x^{(i)} | z^{(i)}=l, \mu, \sigma, \phi) p(z^{(i)}=l | \phi)} \end{aligned}$$

$$\begin{aligned} p(z^{(i)}=j | \phi) &= \phi_j \\ p(z^{(i)}=l | \phi) &= \phi_l \quad \text{proportional} \\ p(x^{(i)} | z^{(i)}=j, \mu, \sigma, \phi) &\propto \exp\left(-\|x^{(i)} - \mu_j\|^2 / \sigma^2\right) \end{aligned}$$

m-step = given  $w_j^{(i)}$ ,  $i=1\dots n$ ,  $j=1\dots K$

$$\phi_j = \frac{1}{n} \sum_{i=1}^n w_j^{(i)} \quad \mu_j = \frac{\sum_{i=1}^n w_j^{(i)} x^{(i)}}{\sum_{i=1}^n w_j^{(i)}}$$

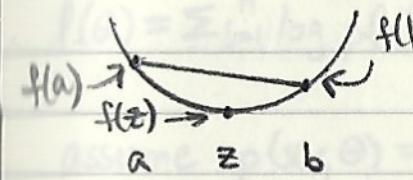
Convex Analysis = a set  $\Omega$  is convex if  $a, b \in \Omega$  the line between  $a, b$  is in  $\Omega$ .



$$\lambda \in [0,1], a, b \in \mathbb{R} \Rightarrow \lambda a + (1-\lambda)b \in \Omega$$

Given function  $f$ , graph of  $f$ ,  $G_f = \{(x, y) : y \geq f(x)\}$

function is convex if its graph is a convex set



$$\lambda f(a) + (1-\lambda)f(b) \in \Omega$$

$$z = \lambda a + (1-\lambda)b$$

$$\lambda f(a) + (1-\lambda)f(b) \geq f(z)$$

for every  $f''(x) > 0$  (twice differentiable)  $\Rightarrow f$  is convex

$$\text{Taylor's theorem: } f(a) = f(z) + f'(z)(a-z) + \underbrace{\frac{1}{2}f''(z_a)(a-z)^2}_{C_a}$$

$$f(b) = f(z) + f'(z)(b-z) + \underbrace{\frac{1}{2}f''(z_b)(b-z)^2}_{C_b}$$

$$\lambda f(a) + (1-\lambda)f(b) = f(z) + f'(z)(\lambda a + (1-\lambda)b - z) + \underbrace{C_a + C_b}_{\text{both } > 0} = 0$$

$$\therefore f(z) < \lambda f(a) + (1-\lambda)f(b)$$

Jensen's inequality:  $\mathbb{E}[f(x)] \geq f(\mathbb{E}[x])$ ,  $x$  is random  
 $f$  is convex

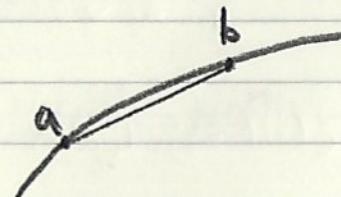
$$\begin{matrix} x \rightarrow a & \lambda \\ & \searrow \\ & b & 1-\lambda \end{matrix}$$

$$\mathbb{E}[f(x)] = \lambda f(a) + (1-\lambda)f(b) \geq f(z)$$

$$\mathbb{E}[x] = \lambda a + (1-\lambda)b = z$$

$f$  is concave if  $-f$  is convex

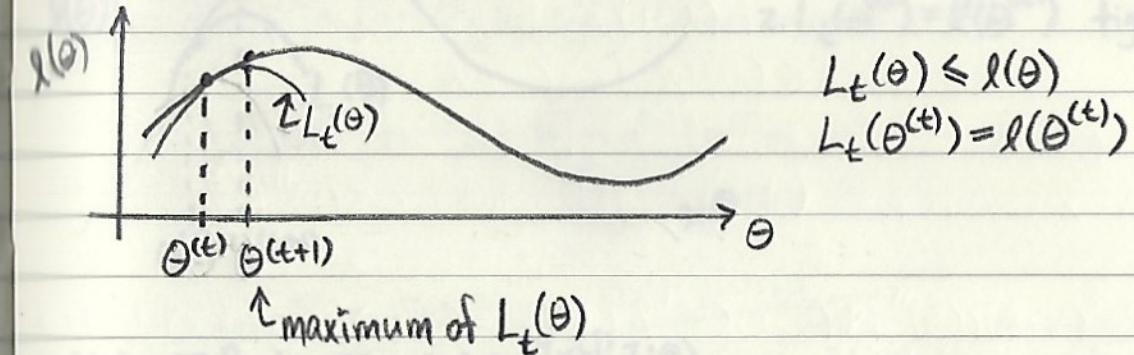
$$\begin{aligned} f &= \log x \\ f'' &= -x^{-2} \end{aligned}$$



$$\mathbb{E}[g(x)] \leq g(\mathbb{E}[x])$$

$$l(\theta) = \sum_{i=1}^n \log p(x^{(i)}; \theta)$$

$$\text{assume } p(x; \theta) = \sum_z p(x, z; \theta)$$



c-step = find  $L_t(\theta^{(t)})$  given  $\theta^{(t)}$

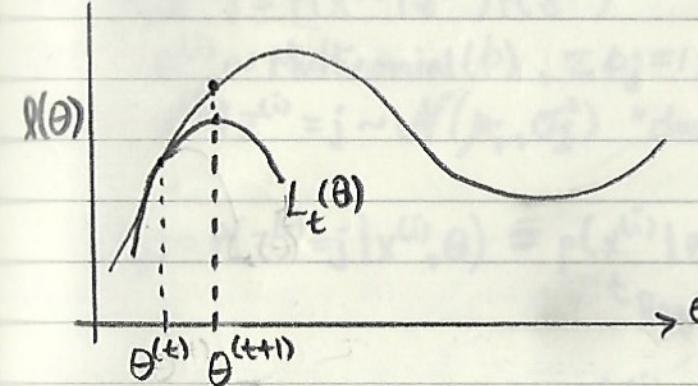
m-step = solve  $\theta^{(t+1)} = \arg \max_{\theta} L_t(\theta^{(t)})$

$$\begin{aligned} \log \sum_z p(x, z; \theta) &= \log \sum_z \frac{Q(z)p(x, z; \theta)}{Q(z)} Q(z) \text{ s.t. } Q(z) \geq 0 \\ &= \log \mathbb{E}_Q \left[ \frac{p(x, z; \theta)}{Q(z)} \right] \leftarrow \text{concave} \end{aligned}$$

$$= \mathbb{E}_Q \left[ \log \frac{p(x, z; \theta)}{Q(z)} \right]$$

$$= \sum_z Q(z) \log \frac{p(x, z; \theta)}{Q(z)}$$

# Machine Learning Expectation Maximization, Factor Analysis



- 1.  $L_t(\theta) \leq l(\theta)$  lower bound
- 2.  $L_t(\theta^{(t)}) = l(\theta^{(t)})$  tight

$$l(\theta) = \sum_{i=1}^n \log \sum_z Q_i(z) \frac{P(x^{(i)}, z; \theta)}{Q_i(z)} \quad \sum_z Q_i(z) = 1, Q_i(z) > 0$$

$$L_t(\theta) = \sum_{i=1}^n \sum_z Q_i(z) \log \frac{P(x^{(i)}, z; \theta)}{Q_i(z)} \leftarrow \text{less than } l(\theta) \forall \theta$$

choose  $Q_i(z)$  - may depend on  $x^{(i)}, \theta$  (do not want to depend on  $z$ )

$$\frac{P(x^{(i)}, z; \theta)}{Q_i(z)} = c \Rightarrow l(\theta) = \sum_{i=1}^n \log \sum_z Q_i(z) \cdot c = \sum_{i=1}^n \log c$$

guarantees  $\nearrow$

$$L_t(\theta) = \sum_{i=1}^n \sum_z Q_i(z) \log c = \sum_{i=1}^n \log c$$

$$Q_i(z) = P(z | x^{(i)}, \theta)$$

e-step: for  $i = 1 \dots n$

$$\text{set } Q_i(z) = P(z | x^{(i)}, \theta^{(t)})$$

$$\text{m-step: } \theta^{(t+1)} = \arg \max L_t(\theta)$$

terminate?  $= l(\theta^{(t+1)}) \geq \max_{\theta} L_t(\theta) \geq L_t(\theta^{(t)}) = l(\theta^{(t)})$  - climbing locally  
not globally optimal

Mixture of Gaussians:

$$P(x^{(i)}, z^{(i)}) = P(x^{(i)} | z^{(i)}) P(z^{(i)})$$

$z^{(i)} \sim \text{Multinomial}(\phi)$ ,  $\sum \phi_j = 1$ ,  $\phi_j > 0$        $z^{(i)}$  "hidden"

$x^{(i)} | z^{(i)} = j \sim \mathcal{N}(\mu_j, \sigma_j^2)$  "cluster means"

$$Q_i(j) = P(z^{(i)}=j | x^{(i)}, \theta) = \frac{p(x^{(i)} | z^{(i)}=j, \theta)}{\text{Bayes' Rule}}$$

$$\max_{\phi, \mu, \Sigma} \sum_{i=1}^n \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}, \theta)}{Q_i(z)}$$

$$\sum_{i,j} w_j^{(i)} \log \frac{1}{\sqrt{2\pi|\Sigma_j|}} \exp\left(-\frac{1}{2}(x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j)\right) \phi_j / w_j^{(i)}$$

$$= \sum_{i,j} w_j^{(i)} \log \frac{1}{\sqrt{2\pi|\Sigma_j|}} = \frac{1}{2} (x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j) + \log \phi_j - \log w_j^{(i)}$$

$$\nabla_{\mu_j} L_t = \sum_i w_j^{(i)} \sum_j (x^{(i)} - \mu_j) = \sum_j (\sum_i w_j^{(i)} (x^{(i)} - \mu_j)) = 0$$

$$\mu_j = \frac{\sum_i w_j^{(i)} x^{(i)}}{\sum_i w_j^{(i)}} \quad \text{↑ full rank because can be inverted}$$

Factor Analysis

$n \ll d$  - many fewer points than dimensions

ex: temperature sensors,  $d = 10,000$ ,  $n \sim 30$

GMM will not work well

$$\mu = \frac{1}{n} \sum x^{(i)} - \text{OK}$$

$$\Sigma = \frac{1}{n} \sum (x^{(i)} - \mu)(x^{(i)} - \mu)^T$$

↑ rank  $= d \leq n < d$  - cannot invert  $\Sigma$  and  $|\Sigma| = 0$

building block 1:  $\Sigma = \sigma^2 I$  - constrain  $\Sigma$ ,  $\Sigma \in \mathbb{R}^{d \times d}$   $\therefore |\Sigma| = \sigma^{2d}$

$$\begin{aligned} \max_{\mu} \sum_{i=1}^n -\log 2\pi |\Sigma|^{\frac{1}{2}} + (x^{(i)} - \mu)^T \Sigma^{-1} (x^{(i)} - \mu) \\ = \max_{z} \sum_{i=1}^n -d \log z + z^{-1} (x^{(i)} - \mu)^T (x^{(i)} - \mu), z = \sigma^2 \end{aligned}$$

$$z = \frac{1}{nd} \sum_{i=1}^n (x^{(i)} - \mu)^T (x^{(i)} - \mu)$$

building block 2:  $\Sigma = \begin{bmatrix} \sigma_1^2 & & \\ & \ddots & \\ & & \sigma_d^2 \end{bmatrix}$  - breaks down into  $d$  1-dimensional problems

$$z_j = \frac{1}{n} \sum_{i=1}^n (x_j^{(i)} - \mu_j)^2$$

model:  $P(x, z) = P(x|z)p(z)$ ,  $z \sim N(0, I) \in \mathbb{R}^s$ ,  $s \ll d$

$$x = \mu + \Lambda z + \epsilon$$

$\Lambda \in \mathbb{R}^{d \times s}$

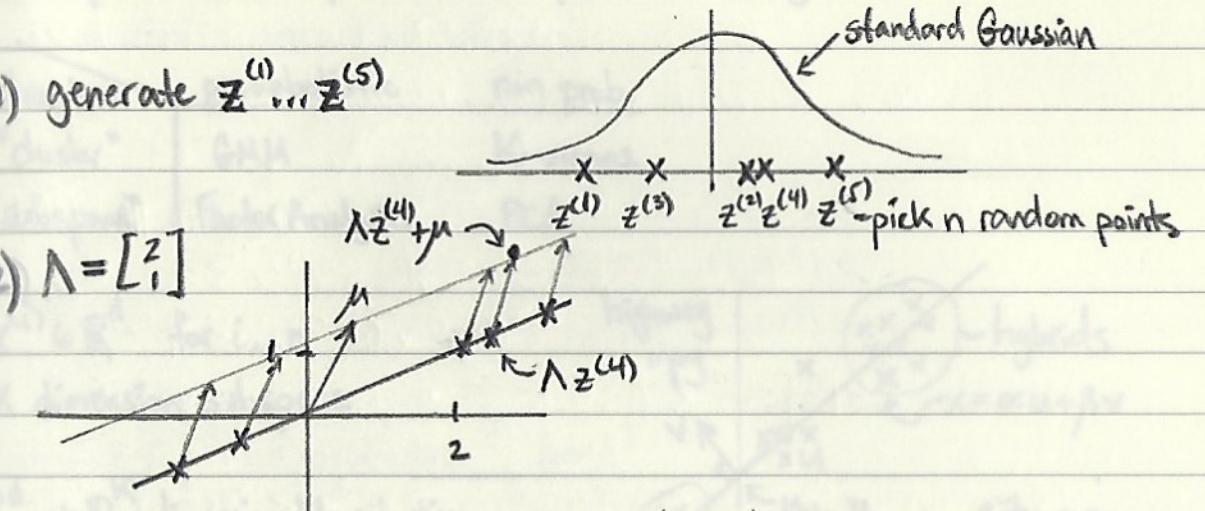
$$\epsilon \sim N(0, \Sigma)$$

$\Sigma$  diagonal

$d=2, s=1, n=5$

i) generate  $z^{(1)} \dots z^{(5)}$

ii)  $\Lambda = [z]$



$$X = \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} \quad \begin{matrix} \uparrow d_1 \\ \uparrow d_2 \end{matrix} \quad \begin{matrix} X_1 \in \mathbb{R}^{d_1} \\ X_2 \in \mathbb{R}^{d_2} \\ X \in \mathbb{R}^{d_1+d_2} \end{matrix}$$

$$\Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \quad \begin{matrix} \uparrow d_1 \\ \uparrow d_2 \end{matrix} \quad \Sigma_{ij} \in \mathbb{R}^{d_i \times d_j}$$

$$X \sim N(\mu, \Sigma)$$

$$\text{fact 1: } p(x_1) = \int_{x_2} p(x_1, x_2) = N(\mu_{11}, \Sigma_{11})$$

$$\text{fact 2: } p(x_1 | x_2) \sim N(\mu_{1|2}, \Sigma_{1|2})$$

$$\begin{aligned} \mu_{1|2} &= \mu_1 + \Sigma_{12}^{-1} \Sigma_{22} (x_2 - \mu_2) \\ \Sigma_{1|2} &= \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \end{aligned}$$

$$\begin{pmatrix} z \\ x \end{pmatrix} \sim N((\mu), \Sigma), \quad \Sigma_{11} = [zz^T] = I$$

$$\begin{aligned} \Sigma_{12} &= \mathbb{E}[z(x-\mu)^T] = \mathbb{E}[z(\lambda z + \epsilon)^T] \\ &= \mathbb{E}[zz^T]\lambda^T + \mathbb{E}[z\epsilon^T] = \lambda^T \end{aligned}$$

$$\Sigma = \begin{bmatrix} I & \lambda^T \\ \lambda & \lambda\lambda^T + \Phi \end{bmatrix}$$

$$\begin{aligned} \Sigma_{22} &= \mathbb{E}[(x-\mu)(x-\mu)^T] = \mathbb{E}[(\lambda z + \epsilon)(\lambda z + \epsilon)^T] \\ &= \mathbb{E}[\lambda zz^T\lambda] + \mathbb{E}[\epsilon\epsilon^T] = \lambda\lambda^T + \Phi \end{aligned}$$

## Machine Learning

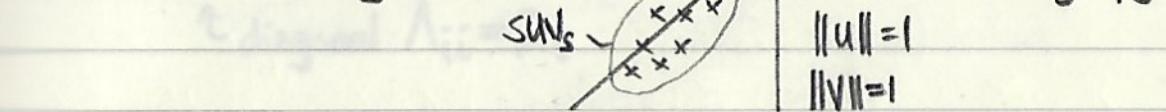
CS229-L14

Principal and Independent Component Analysis

structure	probabilistic	non prob.
"cluster"	GMM	K-means
"subspace"	Factor Analysis	PCA

$x^{(i)} \in \mathbb{R}^d$  for  $i \dots n$   
 $K$  dimension subspace

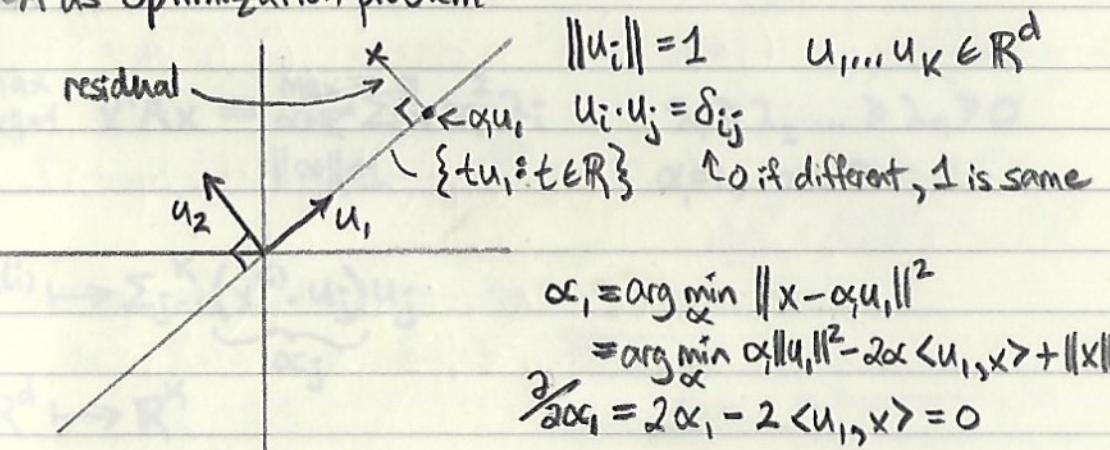
$\mathbb{R}^d \rightarrow \mathbb{R}^K$  dimensionality reduction



preprocessing -  $\mu = \frac{1}{n} \sum_i x^{(i)}$  "centering"  
 $x^{(i)} \mapsto x^{(i)} - \mu$

data should also be of "same scale"

PCA as optimization problem:



$$\min_{\alpha \in \mathbb{R}^K} \|x - \sum_{j=1}^K \alpha_j u_j\|^2 = \min_{\alpha \in \mathbb{R}^K} \left[ \sum_{j=1}^K \alpha_j^2 - 2 \sum_{j=1}^K \alpha_j \langle u_j, x \rangle + \|x\|^2 \right]$$

$$\alpha_j = u_j \cdot x$$

We can find PCA

- 1) maximize amount in subspace
- 2) minimize residual

$$\max_{\substack{u \in \mathbb{R}^d \\ \|u\|=1}} \frac{1}{n} \sum_{i=1}^n (u \cdot x^{(i)})^2 = \frac{1}{n} \sum_{i=1}^n (u \cdot x^{(i)})(x^{(i)} \cdot u) \\ = \frac{1}{n} u^T (\sum_{i=1}^n x^{(i)} x^{(i)T}) u \\ = u^T (\frac{1}{n} \sum_{i=1}^n x^{(i)} x^{(i)T}) u$$

$\uparrow$  covariance matrix

Let  $A \in \mathbb{R}^{d \times d}$  be symmetric,  $A = A^T$

$$A = U \Lambda U^T \quad U U^T = U^T U = I$$

$\uparrow$  diagonal  $\Lambda_{ii} = \lambda_i$

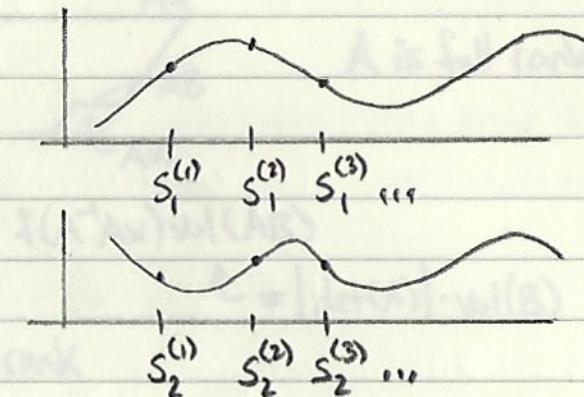
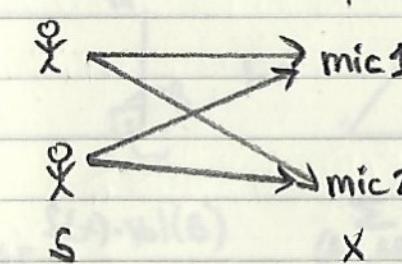
$$Ax = U \Lambda U^T x \quad x = \sum \alpha_i u_i \\ = U \Lambda \sum_{i=1}^d \alpha_i e_i \\ = U \sum \alpha_i \lambda_i e_i \quad \uparrow \text{standard unit vectors} \\ = \sum_{i=1}^d \alpha_i \lambda_i u_i$$

$$\max_{\|x\|=1} x^T A x \equiv \max_{\|\alpha\|=1} \sum_{i=1}^n \alpha_i^2 \lambda_i \quad \lambda_1 \geq \lambda_2 \dots \geq \lambda_n \geq 0 \\ \alpha_i = 1, \text{ or } 0$$

$$x^{(i)} \mapsto \sum_{j=1}^k \underbrace{(\underbrace{x^{(i)}, u_j}_{\alpha_j})}_{\alpha_j} u_j$$

$$\mathbb{R}^d \mapsto \mathbb{R}^k$$

Cocktail party problem:



$$\text{observed: } x_j^{(t)} = \alpha_{j1} s_1^{(t)} + \alpha_{j2} s_2^{(t)}$$

$$X^{(t)} = AS^{(t)}, A \in \mathbb{R}^{d \times d}, S \in \mathbb{R}^d \leftarrow \text{matrix form}$$

$\uparrow$  static

$$\text{given: } X^{(1)}, \dots, X^{(n)} \in \mathbb{R}^d \quad (\text{d mics})$$

$$\text{find: } S^{(1)}, \dots, S^{(d)} \in \mathbb{R}^d \quad (\text{d speakers})$$

$A \in \mathbb{R}^{d \times d}$  - mixing matrix,  $W = A^{-1}$  - unmixing matrix

$$S^{(t)} = A^{-1} X^{(t)}$$

$$W = \begin{bmatrix} W_1^T \\ \vdots \\ W_d^T \end{bmatrix} \quad S_j^{(t)} = W_j \cdot X^{(t)}$$

$S$  cannot be Gaussian -  $S^{(i)} \sim N(0, I) \therefore X^{(i)} \sim N(0, AA^T)$

$$AA^T = AIA^T = AA^T$$

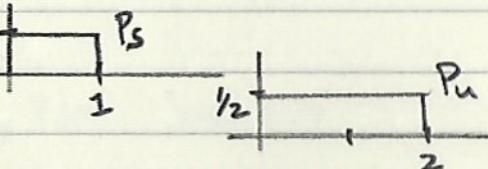
density under linear transform =

$$\text{ex: } S \sim \text{uniform}(0, 1)$$

$$U = 2S$$

$$P_U(\frac{1}{2}) \neq P_S(x)$$

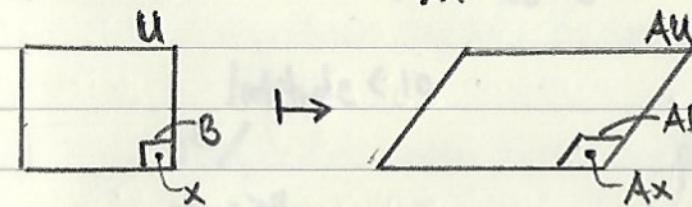
$$P_U(x) = P_S(A^{-1}x) \cdot |\det(A^{-1})|$$



CS229 = L14

$$\int_U f(x) dx$$

$$\int_{A^T} f(A^{-1}y) dy = |\det(A)| \cdot \int_U f(x) dx$$



A is full rank

$$\sum_{(x, B) \in P} f(x) \cdot \text{vol}(B)$$

$$\sum_{(Ax, AB) \in \hat{P}} f(A^{-1}Ax) \text{vol}(AB)$$

$$\uparrow = |\det(A)| \cdot \text{vol}(B)$$

$$\det(A^{-1}) = \frac{1}{\det(A)} \text{ with } A \text{ full rank}$$

$$P_S(S) = \prod_{i=1}^d P_S(S_i)$$

$$p(x) = \prod_{j=1}^d P_S(w_j \cdot x) / |\det(W)|$$

$$P_S(x) \propto g'(x) \text{ for } g(x) = (1 + e^k)^{-1}$$

$$\max_{\Theta} \log P(x, \Theta)$$

## Machine Learning Decision Trees, Boosting, Bagging

CS229: S7

latitude &lt; 10

yes /

time &gt; 4

yes / no \

time &lt; 8 (-)

yes / no \

(+ ) (- )

prediction - drop down tree  
until reach node

impurity as a loss function:

- misclassification loss - # examples misclassified with majority prediction

$$L_{\text{misclass}}(\hat{p}) = 1 - \max(\hat{p}, 1 - \hat{p})$$

- cross-entropy loss

$$L_{\text{cross}}(\hat{p}) = -\hat{p} \log \hat{p} - (1 - \hat{p}) \log (1 - \hat{p})$$

when to split leaf node? - decrease average cardinality-weighted loss in children

$$\frac{|R_1| L(R_1) + |R_2| L(R_2)}{|R_1| + |R_2|}$$

handles categorical variables (e.g. North, South)

regularization: max leaf size  
max depth

advantages = interpretable  
easily to visualize  
categorical data

disadvantages = high variance  
optimal NP-complete

7/15/19

selected models, 7/14/19

Prof. Christopher R.

11/8/19

#2: Pass

prior knowledge

correct, confident, want nice G

01 &gt; 3rd digit

not much job - not strong  
when they know

\approx

Kent

/n \approx

- R&gt;0

/n \approx

-

+

not bad, not a good  
but better than others -  
adaboost, phrasing

- (q-1, q) when -1 = (q)

(q-1)q(q-1) - q(q-1) = (q)q(q-1) -

and helping - standard approach - when first step ok -  
multiple i

(2).1.91 + (2).1.91

1.91 + 1.91

(Adaboost) adding together without

use full sum = with values

high sum

sum over all = regular basis

selected weights = regular basis

with values of loss

other losses too

CS229: SF

bias - classifier not expressive enough

variance - final hypothesis depends on sampling of training data

bagging - resample data many times to reduce variance  
average outputs of trained models

e.g. random forest

boosting - design a sequential process where weak learners  
try solving problem. used with low-bias models.

additive modeling -

input: training set  $\{(x^{(i)}, y^{(i)})\}_{i=1}^m$  and distribution  $p^{(1)}, \dots, p^{(m)}$   
with  $\sum_{i=1}^m p^{(i)} = 1$  and  $p^{(i)} \geq 0$ weak classifier:  $\phi_j: \mathbb{R}^n \rightarrow \{-1, 1\}$  s.t.  $\sum_{i=1}^m p^{(i)} \mathbb{1}\{y^{(i)} \neq \phi_j(x^{(i)})\} \leq \frac{1}{2} - \gamma$ 

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \exp(-y^{(i)} \theta^T \phi(x^{(i)}))$$

$$\frac{\partial J}{\partial \theta} = \sum w^+ e^{-\theta^T \phi} + \sum w^- e^{\theta^T \phi}$$

Adaboost - additive modeling with exponential loss

Gradient boost - use gradients to create weak classifiers

## Machine Learning Weak Supervision

CS229-L15

active learning = select points to label more intelligently

semi-supervised learning = use unlabeled data as well

transfer learning = transfer from one dataset to another

weak supervision = label data in cheaper, higher-level ways

- ↳ programmatic supervision - heuristics to label training data

- ↳ heuristic supervision

- distant supervision

- data augmentation

- other models

### ML Applications - model + data + hardware

↑ comes from dirty, messy processes

label quality and quantity much more important than model choice

manual labels - slow, expensive, static (relabel if schema changes)

programmatic labels - fast, cheap, dynamic ← lower quality

Snorkel = formalizing programmatic labeling

↳ replace ad hoc weak supervision

problem = noisy sources conflict and are correlated

- ↳ model and combine noisy labels into probabilities

- ↳ use probabilistic labels to train a model

key idea = observe overlapping judgements on many points

given:  $x^{(1)}, \dots, x^{(n)} \in \mathbb{R}^d$

$\lambda_1, \dots, \lambda_m, \lambda_i: \mathbb{R}^d \rightarrow \{-1, 1\}$  ← labeling functions

find:  $p(y^{(i)} = 1 | \bar{\lambda}_i, x^{(i)})$  ← conditionally independent  
↑ all labeling functions

model = each labeler has accuracy  $\alpha_i$  (hidden)

with probability  $\alpha_i$ ,  $\lambda_i(x) = y$  ← correct

" " "  $1 - \alpha_i$ ,  $\lambda_i(x) = -y$  ← incorrect

$$\mathbb{E}[\lambda_i; y] = \alpha_i \cdot 1 + (1 - \alpha_i) \cdot (-1) = 2\alpha_i - 1 = \alpha_i$$

↑ agree ↑ disagree

$$\mathbb{E}[\lambda_i \lambda_j] = \mathbb{E}[\lambda_i Y \lambda_j Y] = \begin{cases} \alpha_i \alpha_j & i \neq j \\ 1 & \text{o.w.} \end{cases}$$

$$M_{ij} = \mathbb{E}[\lambda_i \lambda_j] \leftarrow \text{observe}$$

$$M_{ij} = \sum_{k=1}^n \lambda_j(x^{(i)}) \lambda_k(x^{(j)})$$

$$i, j, k \text{ distinct indexes} - \frac{M_{ij} M_{ik}}{M_{kk}} = \frac{(\alpha_i \alpha_j)(\alpha_i \alpha_k)}{\alpha_i \alpha_k} = \alpha_j^2$$

↑ 3 distinct labelers "vote"

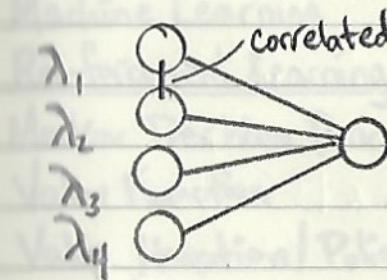
Knowing sign of  $\alpha_i$  enables knowing sign for all  $\alpha_j$   
 $\alpha_i = 0$  means accuracy is  $1/2$  ← random

gradient descent:

$$\min_a \|M - (I + \begin{bmatrix} \alpha_1^2 & & \\ & \ddots & \\ & & \alpha_n^2 \end{bmatrix} + aa^T)\|_F^2$$

212:95529

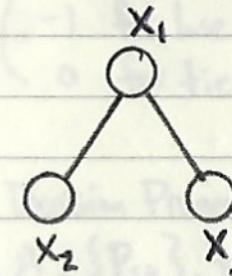
CS229-L15



$$\mathbb{E}[\lambda_i \lambda_j | Y] = \mathbb{E}[\lambda_i | Y] \mathbb{E}[\lambda_j | Y] \quad (i, j) \notin E$$

← Graphical model

use of inverse co-variance:



$$\begin{aligned} x_1 &\sim N(0, 1) & \mathbb{E}[x_1] &= 0 \\ x_2 &\sim N(x_1, 1) & x_2 &= x_1 + \epsilon_2 \\ x_3 &\sim N(x_1, 1) & x_3 &= x_1 + \epsilon_3 \xrightarrow{\text{---}} N(0, 1) \end{aligned}$$

$$\mathbb{E}[x_1^2] = 1$$

$$\mathbb{E}[x_2^2] = \mathbb{E}[(x_1 + \epsilon_2)^2] = 2$$

$$\mathbb{E}[x_1 x_2] = \mathbb{E}[x_1^2 + x_1 \epsilon_2] = 1$$

$$\mathbb{E}[x_2 x_3] = \mathbb{E}[(x_1 + \epsilon_2)(x_1 + \epsilon_3)] = 1$$

$$\Sigma = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 1 \\ 1 & 1 & 2 \end{bmatrix}, \Sigma^{-1} = \begin{bmatrix} 3 & -1 & -1 \\ -1 & 1 & 0 \\ -1 & 0 & 1 \end{bmatrix}$$

zeros where no ↑  
edges present

We say a  $p: \mathbb{R}^d \rightarrow [0, 1]$  agrees with graph

$$G = (\{1 \dots d\}, E) \text{ if } p(x) = \prod_{\substack{u, v \in V \\ (u, v) \in E}} P_{u,v}(x_u, x_v) \cdot \prod_{v \in V} P_v(x_v)$$

## Machine Learning Reinforcement Learning

Markov Decision Processes - MDPs

Value Function

Value Iteration / Policy Iteration

$$R(s) = \begin{cases} +1 & \text{for win} \\ -1 & \text{for lose} \\ 0 & \text{for tie} \end{cases}$$

Reward as function of state

Markov Decision Process:

$$(S, A, \{P_{SA}\}, \gamma, R)$$

$S$ : set of states

$A$ : set of actions

$P_{SA}$ : state transition probability function,  $\sum_s P_{sa}(s') = 1$

$\gamma$ : discount factor  $\in (0, 1]$

$R$ : rewards

$\rightarrow$	$\rightarrow$	$\rightarrow$	$+1 \downarrow$	terminal states
$\uparrow$		$\uparrow$	$-1 \downarrow$	
$\frac{\oplus}{\ominus} \uparrow$	$\leftarrow$	$\leftarrow$	$\leftarrow$	
$\rightarrow$	$\rightarrow$	$\rightarrow$	$+1 \downarrow$	
$\uparrow$		$\uparrow$	$-1 \downarrow$	
$\frac{\oplus}{\ominus} \uparrow$	$\leftarrow$	$\leftarrow$	$\leftarrow$	

optimal policy

II states  
Actions:  $\{N, S, E, W\}$

$$R((4, 3)) = +1$$

$$R((4, 2)) = -1$$

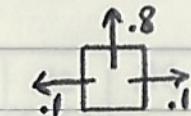
$$R(s) = -0.02 \text{ for all other } s$$

↑ incentive to get to target

$$P_{(3,1)N} = ((3, 2)) = 0.8$$

$$P_{(3,1)N} = ((4, 1)) = 0.1$$

$$P_{(3,1)N} = ((2, 1)) = 0.1$$



start  $s_0$

choose action  $a_0$

get to  $s_1 \sim P_{s_0, a_0}(s_1)$

choose action  $a_1$

get to  $s_2 \sim P_{s_1, a_1}(s_2)$

total payoff:

$$R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots \quad \text{discounting future rewards}$$

↳ required for convergence

goal = choose actions over time to maximize reward

$$\mathbb{E}[R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots]$$

policy  $\pi: S \mapsto A$

controller

define  $V^\pi: S \mapsto \mathbb{R}$  for policy  $\pi$

$V^\pi(s)$  is expected total payoff for starting in  $s$

and executing policy  $\pi$  - value function

$$V^\pi(s) = \mathbb{E}[R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots | \pi, s_0 = s]$$

$\pi$	$\rightarrow$	$\rightarrow$	$\rightarrow$	$+1$	$\sqrt{\pi}$	.52	.73	.77	+1
	$\downarrow$		$\rightarrow$	-1		-.90		-.82	-1
	$\rightarrow$	$\rightarrow$	$\uparrow$	$\uparrow$		-.88	-.87	-.85	-1.00

bellman's equation:

$$V^\pi(s) = R(s) + \gamma \sum_{s'} P_{s, a}(s') V^\pi(s')$$

## Bellman's Equations:

$$V^\pi(s) = R(s) + \gamma \sum_{s'} P_{s,\pi(s)}(s') V^\pi(s')$$

future reward

$$\begin{aligned} V^\pi(s) &= \mathbb{E}[R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots | \pi, s_0 = s] \\ &= R(s) + \gamma \mathbb{E}[R(s_1) + \gamma R(s_2) + \dots | \pi, s_0 = s] \\ &= R(s) + \gamma \mathbb{E}[V^\pi(s_1) | \pi, s_0 = s] \\ &= R(s) + \gamma \sum_{s'} P_{s,\pi(s)}(s') V^\pi(s') \end{aligned}$$

$$V^\pi(s) = R(s) + \gamma \sum_{s'} P_{s,\pi(s)}(s') V^\pi(s')$$

$$\begin{aligned} V^\pi((3,1)) &= R((3,1)) + \gamma (0.8 V^\pi((3,2)) + \\ &\quad 0.1 V^\pi((4,1)) + \\ &\quad 0.1 V^\pi((2,1))) \end{aligned}$$

$$\begin{bmatrix} V^\pi((1,1)) \\ V^\pi((2,1)) \\ \vdots \\ V^\pi((4,1)) \end{bmatrix} \in \mathbb{R}^n$$

$V^*$  : optimum value function  
 $V^*(s) = \max_{\pi} V^\pi(s)$

Bellman's equation :

$$V^*(s) = R(s) + \max_{a \in A} \gamma \sum_{s'} P_{sa}(s') V^*(s')$$

CS221 = L16

$$\pi^*(s) = \arg \max_{a \in A} \sum_{s'} P_{sa}(s') V^*(s') : \text{optimum policy function}$$

$$V^*(s) = V^{\pi^*}(s) \geq V^\pi(s)$$

Value iteration =  $V$  converges to  $V^*$  (gets arbitrarily close)

initialize  $V(s) = 0 \forall s$  (for all  $s$ )

$\forall s$ , update:

$$V(s) := R(s) + \max_a \gamma \sum_{s'} P_{sa}(s') V(s')$$

Synchronous update = update all values of  $V$  simultaneously

asynchronous update = update one  $V$  at a time

Bellman backup operator:  $V := B(V)$

.86	.90	.93	+1
.81		.69	-1
.78	.75	.71	.49

Policy iteration = converges to  $\pi^*$  (get to  $\pi^*$  but requires more computation of Bellman's equation)

repeat

set  $V = V''$  (solve Bellman's equations)

$$\text{set } \pi(s) = \arg \max_{a \in A} \sum_{s'} P_{sa}(s') V(s')$$

## Linear Algebra Review

$$c = AB \in \mathbb{R}^{m \times p}, C_{ij} = \sum_{k=1}^n A_{ik} B_{kj}$$

$$x^T y \in \mathbb{R} = \sum_{i=1}^n x_i y_i$$

$$(AB)C = A(BC), A(B+C) = AB + AC$$

$$AI = A = IA$$

$$(AT)^T = A, (AB)^T = B^T A^T, (A+B)^T = AT + BT$$

Symmetric:  $A = A^T, A \in \mathbb{R}^{n \times n}$

for any  $A \in \mathbb{R}^{n \times n}$ ,  $A + AT$  is symmetric  
 $\hookrightarrow A = \frac{1}{2}(A+AT) + \frac{1}{2}(A-AT)$

$$\text{tr}(A) = \sum_{i=1}^n A_{ii}$$

$$A \in \mathbb{R}^{n \times n}, \text{tr} A = \text{tr} AT$$

$$A, B \in \mathbb{R}^{n \times n}, \text{tr}(A+B) = \text{tr} A + \text{tr} B$$

$$\text{if } AB \text{ is square, } \text{tr} AB = \text{tr} BA$$

$$\text{if } ABC \text{ is square, } \text{tr} ABC = \text{tr} BCA = \text{tr} CAB$$

$$\nabla_{AT} f(A) = (\nabla_A f(AT))^T$$

$$\nabla_A \text{tr} AB = BT$$

$$\nabla_A \text{tr} ABAT^T = CAB + CT ABT$$

$$\nabla_A |A| = |A|(A^{-1})^T - \text{non-singular square } A$$

$$\|x\|_2 = \sqrt{\sum_{i=1}^n x_i^2}, \|x\|_2^2 = x^T x, \|x\|_1 = \sum_{i=1}^n |x_i|$$

$$A \in \mathbb{R}^{m \times n} - \text{col rank} = \text{row rank} = \text{rank}$$

$$\text{rank}(A) \leq \min(m, n), \text{rank}(A) = \min(m, n) \text{ full rank}$$

$$\text{rank}(A) = \text{rank}(AT)$$

$$\text{rank}(AB) = \min(\text{rank}(A), \text{rank}(B))$$

$$A^{-1}A = I = AA^{-1}, A \in \mathbb{R}^{n \times n}$$

$$(A^{-1})^{-1} = A, (AB)^{-1} = B^{-1}A^{-1}, (A^{-1})^T = (AT)^{-1} = A^T$$

$$x^T y = 0 - \text{orthogonal}$$

$U \in \mathbb{R}^{n \times n}$  - orthogonal if all columns orthogonal and normalized

$$U^T U = I = UU^T, \|Ux\|_2 = \|x\|_2 \quad x \in \mathbb{R}^n$$

span of  $\{x_1, \dots, x_n\}$  = set of all vectors that are linear comb.

$$\nabla_\theta J(\theta) = x^T \theta - x^T y = 0$$

$$\theta = (x^T x)^{-1} x^T y$$

Probabilistic Interpretation

$$y^{(i)} = \theta^T x^{(i)} + \epsilon^{(i)}, \epsilon^{(i)} \sim \mathcal{N}(0, \sigma^2)$$

$$P(y^{(i)} | x^{(i)}; \theta) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)$$

$$L(\theta) = \prod_{i=1}^n P(y^{(i)} | x^{(i)}; \theta)$$

$$l(\theta) = \log L(\theta) = n \log \frac{1}{\sigma\sqrt{2\pi}} - \frac{1}{\sigma^2} \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2$$

Locally-weighted linear regression

$$\text{fit } \theta \text{ to } \sum_i w^{(i)} (y^{(i)} - \theta^T x^{(i)})^2$$

$$w^{(i)} = \exp\left(-\frac{1}{2\tau^2} (x^{(i)} - \bar{x})^2\right)$$

$$\text{tr } A = \sum_{i=1}^n \lambda_i, |A| = \prod_{i=1}^n \lambda_i$$

$$\text{Hessian } \nabla_x^2 f(x) \in \mathbb{R}^{n \times n} = \begin{bmatrix} \frac{\partial^2 f(x)}{\partial x_1^2} & \frac{\partial^2 f(x)}{\partial x_1 \partial x_2} & \dots \\ \vdots & \ddots & \\ \frac{\partial^2 f(x)}{\partial x_n \partial x_1} & \dots & \frac{\partial^2 f(x)}{\partial x_n^2} \end{bmatrix}$$

always symmetric

## Probability Theory

$$\mathbb{E}[g(x)] \triangleq \sum_x g(x) p(x)$$

$$\mathbb{E}[a] = a \text{ for constant } a$$

$$\mathbb{E}[af(x)] = a \mathbb{E}[f(x)]$$

$$\mathbb{E}[f(x) + g(x)] = \mathbb{E}[f(x)] + \mathbb{E}[g(x)]$$

$$\mathbb{E}[\mathbb{1}_{\{X=k\}}] = P(X=k)$$

$$\text{Var}[x] \triangleq \mathbb{E}[(x - \mathbb{E}[x])^2] \geq 0$$

$$= \mathbb{E}[x^2] - \mathbb{E}[x]^2$$

$$\text{Var}(x_1 + \dots + x_n) = \text{Var}(x_1) + \dots + \text{Var}(x_n)$$

$x_i$  independent

$$\text{Var}[a] = 0, \text{Var}[af(x)] = a^2 \text{Var}[f(x)]$$

$$\text{Bernoulli: } \begin{array}{ll} \text{mean} & p \\ \text{var} & p(1-p) \end{array}$$

$$\text{Gaussian: } \begin{array}{ll} \mu & \text{mean} \\ \sigma^2 & \text{var} \end{array}$$

$$\text{Poisson: } \begin{array}{ll} \lambda & \text{mean} \\ \lambda & \text{var} \end{array}$$

$$\text{Cov}[x, y] \triangleq \mathbb{E}[(x - \mathbb{E}[x])(y - \mathbb{E}[y])] = \mathbb{E}[xy] - \mathbb{E}[x]\mathbb{E}[y]$$

$$\text{Var}[x+y] = \text{Var}[x] + \text{Var}[y] + 2\text{Cov}[x, y]$$

$$\text{Var}[Ax] = A \Sigma A^T, A \in \mathbb{R}^{m \times n}$$

$$\text{Var}[ax] = a^2 \text{Var}[x], a \in \mathbb{R}^n$$

## LMS Algo

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n (h_\theta(x^{(i)}) - y^{(i)})^2$$

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

$$\frac{\partial}{\partial \theta_j} J(\theta) = \frac{\partial}{\partial \theta_j} \frac{1}{2} (h_\theta(x) - y)^2$$

$$= (h_\theta(x) - y) \frac{\partial}{\partial \theta_j} (h_\theta(x) - y)$$

$$= (h_\theta(x) - y) x_j$$

$$J(\theta) = \frac{1}{2} (x \theta - y)^T (x \theta - y)$$

$$\theta = (x^T x)^{-1} x^T y$$

## Probabilistic Interpretation

$$y^{(i)} = \theta^T x^{(i)} + \epsilon^{(i)}, \epsilon^{(i)} \sim \mathcal{N}(0, \sigma^2)$$

$$P(y^{(i)} | x^{(i)}; \theta) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)$$

$$L(\theta) = \prod_{i=1}^n P(y^{(i)} | x^{(i)}; \theta)$$

$$l(\theta) = \log L(\theta) = n \log \frac{1}{\sigma\sqrt{2\pi}} - \frac{1}{\sigma^2} \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2$$

Locally-weighted linear regression

$$\text{fit } \theta \text{ to } \sum_i w^{(i)} (y^{(i)} - \theta^T x^{(i)})^2$$

$$w^{(i)} = \exp\left(-\frac{1}{2\tau^2} (x^{(i)} - \bar{x})^2\right)$$

## Logistic Regression

$$h_\theta(x) = g(\theta^T x) = \frac{1}{1+e^{-\theta^T x}}$$

$$g(z) = \frac{1}{1+e^{-z}}, g'(z) = g(z)(1-g(z))$$

$$P(y=1 | x; \theta) = h_\theta(x)$$

$$P(y=0 | x; \theta) = 1 - h_\theta(x)$$

$$p(y|x; \theta) = (h_\theta(x))^y (1-h_\theta(x))^{1-y}$$

$$l(\theta) = \sum_{i=1}^n y^{(i)} \log h_\theta(x^{(i)}) + (1-y^{(i)}) \log (1-h_\theta(x^{(i)}))$$

$$\frac{\partial}{\partial \theta_j} l(\theta) = (y - h_\theta(x)) x_j$$

## Newton's Method

$$\theta := \theta - \frac{f(\theta)}{f'(\theta)} = \theta - \frac{\ell'(\theta)}{\ell''(\theta)}$$

$$:= \theta - H^{-1} \nabla_\theta \ell(\theta), H_{ij} = \frac{\partial^2 \ell(\theta)}{\partial \theta_i \partial \theta_j}$$

## Generalized Linear Models

$$p(y; \eta) = b(y) \exp(\eta^T T(y) - a(\eta))$$

$$\text{Bernoulli: } p(y; \phi) = \phi^y (1-\phi)^{1-y}$$

$$= \exp(y \log \phi + (1-y) \log (1-\phi))$$

$$= \exp((\log \frac{\phi}{1-\phi}) y + \log(1-\phi))$$

$$T(y) = y, a(\eta) = -\log(1-\phi), b(y) = 1$$

$$\eta = \log \frac{\phi}{1-\phi} = \log(1+e^\eta), \phi = \frac{1}{1+e^{-\eta}}$$

$$\text{Gaussian: } \sigma^2 = 1$$

$$p(y|\mu) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} (y - \mu)^2\right)$$

$$= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} y^2\right) \cdot \exp(\mu y - \frac{1}{2} \mu^2)$$

$$\eta = \mu, T(y) = y, a(\eta) = \frac{\mu^2}{2} = \frac{\eta^2}{2}$$

$$b(y) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} y^2\right)$$

Assumptions =  $y | x; \theta \sim \text{Exp Family}(\eta)$

2) predict  $E[T(y)]$  given  $x$

$$\text{if } T(y) = y, h(x) = E[y|x]$$

$$3) \eta = \theta^T x$$

ordinary least squares =

$$h(x) = E[y|x] = y = \eta = \theta^T x$$

$$y|x; \theta \sim \mathcal{N}(\mu, \sigma^2)$$

logistic regression =

$$h(x) = E[y|x; \theta] = \phi = \frac{1}{1+e^{-\eta}} = \frac{1}{1+e^{-\theta^T x}}$$

$$y|x; \theta \sim \text{Bernoulli}(\phi)$$

Canonical response  $g(\eta) = E[T(y); \eta]$

canonical link  $g^{-1}$

Properties: max likelihood w.r.t.  $\eta$  concave

$$E[y; \eta] = \frac{\partial \eta}{\partial \eta} a(\eta)$$

$$\text{Var}[y; \eta] = \frac{\partial^2 \eta}{\partial \eta^2} a(\eta)$$

$$\text{Poisson: } \frac{\eta}{\text{log } \lambda} = \frac{a(\eta)}{\eta^2}$$

$$\text{Geometric: } \log(1-\phi) / \log(\frac{e^{\eta}}{1-e^{\eta}}) = \frac{a(\eta)}{\eta}$$

## Softmax regression with GLM:

$$(T(y))_i = \mathbb{I}\{y=i\} - \text{one hot } T(y) \in \mathbb{R}^{K-1}, T(k) = \vec{0}$$

$$\mathbb{E}[T(y)_i] = P(y=i) = \phi_i \\ P(y; \phi) = \phi_1^{\mathbb{I}\{y=1\}} \phi_2^{\mathbb{I}\{y=2\}} \dots \phi_K^{\mathbb{I}\{y=K\}} \\ = \phi_1^{(T(y))_1} \phi_2^{(T(y))_2} \dots \phi_K^{1 - \sum_{i=1}^{K-1} (T(y))_i}$$

$$= \exp((T(y)_1 \log \phi_1 + (T(y))_2 \log \phi_2 + \dots + (1 - \sum_{i=1}^{K-1} (T(y))_i) \log \phi_K)) \\ = \exp((T(y)_1 \log \frac{\phi_1}{\phi_K} + (T(y))_2 \log \frac{\phi_2}{\phi_K} + \dots + (T(y))_{K-1} \log \frac{\phi_{K-1}}{\phi_K} + \log \phi_K))$$

$$\eta = \begin{bmatrix} \log \phi_1 / \phi_K \\ \vdots \\ \log \phi_{K-1} / \phi_K \end{bmatrix}, a(\eta) = -\log(\phi_K), g(\eta_i) = \eta_i = \log \frac{\phi_i}{\phi_K} \\ b(y) = 1 \quad \phi_i = \frac{e^{\eta_i}}{\sum_{j=1}^K e^{\eta_j}}$$

$$\eta_i = \Theta_i^T x, p(y=i|x; \theta) = \phi_i = \frac{e^{\Theta_i^T x}}{\sum_{j=1}^K e^{\Theta_j^T x}}$$

$$\eta_\theta(x) = \mathbb{E}[T(y)|x; \theta] = \mathbb{E}\left[\begin{array}{l} \mathbb{I}\{y=1\} \\ \vdots \\ \mathbb{I}\{y=K\} \end{array} \mid x; \theta\right] = \begin{bmatrix} \phi_1 \\ \vdots \\ \phi_{K-1} \end{bmatrix}, P(y=k|x; \theta) = 1 - \sum_{i=1}^{K-1} \phi_i$$

$$\ell(\theta) = \sum_{i=1}^n \log p(y^{(i)}|x^{(i)}; \theta) = \sum_{i=1}^n \log \prod_{k=1}^K \left( \frac{e^{\Theta_k^T x^{(i)}}}{\sum_{j=1}^K e^{\Theta_j^T x^{(i)}}} \right)$$

## Multivariate Gaussian

$$p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)\right)$$

$$\Sigma \text{ is PD, so } \Sigma^{-1} \text{ is PD} \quad \Sigma_{ij} = \text{Cov}[x_i, x_j]$$

$$\Sigma = \mathbb{E}[(x - \mu)(x - \mu)^T] = \mathbb{E}[xx^T] - \mu\mu^T = \text{Cov}[x_i, x_j]$$

$\Sigma$  is PSD for any random  $X$

$$\text{univariate } \sigma^2 = \mathbb{E}[(x - \mu)^2] = \mathbb{E}[x^2] - \mu^2, \mu = \mathbb{E}[x]$$

## Generative Learning Algorithms

discriminative - learn  $p(y|x)$  or  $h_g(y)$  directly

generative - learn  $p(x|y) \neq p(y)$

$$\arg \max_y p(y|x) = \arg \max_y p(x|y)p(y)$$

$$X \sim N(\mu, \Sigma), \mathbb{E}[x] = \int_x p(x; \mu, \Sigma) dx = \mu$$

$$\text{Cov}(z) = \mathbb{E}[zz^T] - \mathbb{E}[z](\mathbb{E}[z])^T$$

$$\text{Cov}(x) = \Sigma$$

$$y \sim \text{Bernoulli}(\phi), x|y=0 \sim N(\mu_0, \Sigma) \\ x|y=1 \sim N(\mu_1, \Sigma)$$

$$p(y) = \phi^y (1-\phi)^{1-y}$$

$$p(x|y=0) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (x - \mu_0)^T \Sigma^{-1} (x - \mu_0)\right)$$

$$\ell(\phi, \mu_0, \mu_1, \Sigma) = \log \prod_{i=1}^n p(x^{(i)}|y^{(i)}; \mu_0, \mu_1, \Sigma) p(y^{(i)}; \phi)$$

$$\text{maximize } \ell: \phi = \frac{1}{n} \sum_{i=1}^n \mathbb{I}\{y^{(i)}=1\}$$

$$\mu_0 = \frac{1}{n} \sum_{i=1}^n \mathbb{I}\{y^{(i)}=0\} \\ \Sigma = \frac{1}{n} \sum_{i=1}^n (x^{(i)} - \mu_0)(x^{(i)} - \mu_0)^T$$

decision boundary with common  $\Sigma$  orthogonal to line connecting  $\mu_0$ 's

$$p(y=1|x; \phi, \Sigma, \mu_0, \mu_1) = \frac{1}{1 + \exp(-\Theta^T x)} \quad (x_0=1)$$

$\Theta$  approximate function of  $\phi, \mu_0, \mu_1, \Sigma$

## Naive Bayes

$$x \in \mathbb{S}^d, d=50,000, x_i \in \mathbb{S}^{d_i}$$

$$p(x_1, \dots, x_d | y) = \prod_{i=1}^d p(x_i | y)$$

$$\mathcal{L}(\phi_y, \phi_{j|y=0}, \phi_{j|y=1}) = \prod_{i=1}^n p(x^{(i)}|y^{(i)})$$

$$\text{maximize} = \phi_{j|y=0} = \frac{\sum_{i=1}^n \mathbb{I}\{x_i^{(i)}=1 \wedge y^{(i)}=0\} + 1}{\sum_{i=1}^n \mathbb{I}\{y^{(i)}=0\}} = \frac{0.3}{2} + 1$$

$$\phi_y = \frac{1}{n} \sum_{i=1}^n \mathbb{I}\{y^{(i)}=1\} \quad \uparrow \text{possible values of } y$$

Multinomial =  $\{x^{(i)}, y^{(i)}\}; i=1 \dots n$  where  $x^{(i)} = (x_1^{(i)}, x_2^{(i)}, \dots, x_{d_i}^{(i)})$ ,  $d = \# \text{words in } x^{(i)}$

$$\mathcal{L}(\phi_y, \phi_{k|y=0}, \phi_{k|y=1}) = \prod_{i=1}^n p(x^{(i)}, y^{(i)})$$

$$= \prod_{i=1}^n \left( \prod_{j=1}^{d_i} p(x_j^{(i)}|y^{(i)}; \phi_{k|y=0}, \phi_{k|y=1}) \right) p(y^{(i)}; \phi_y)$$

$$\text{maximize} = \phi_{k|y=0} = \frac{\sum_{i=1}^n \sum_{j=1}^{d_i} \mathbb{I}\{x_j^{(i)}=1 \wedge y^{(i)}=0\} + K}{\sum_{i=1}^n \sum_{j=1}^{d_i} \mathbb{I}\{y^{(i)}=0\} + |V|} = \frac{0.3}{d_i} + 1$$

## Kernel Methods

$\phi: \mathbb{R}^d \mapsto \mathbb{R}^P$  may be high or infinite dimension

$$\phi := \Theta + \alpha \sum_{i=1}^n (y^{(i)} - \Theta^T \phi(x^{(i)})) \phi(x^{(i)})$$

$$\Theta = \sum_{i=1}^n \beta_i \phi(x^{(i)}), \Theta \in \mathbb{R}^P, \beta \in \mathbb{R}^n$$

$$\Theta := \sum_{i=1}^n (\beta_i + \alpha(y^{(i)} - \Theta^T \phi(x^{(i)}))) \phi(x^{(i)})$$

$$\phi_i := \phi_i + \alpha(y^{(i)} - \Theta^T \phi(x^{(i)})) \quad K_{ij} = \phi(x^{(i)})^T \phi(x^{(j)}) = \phi(x^{(i)})^T \phi(x^{(i)}) = K_{ii}$$

$$:= \phi_i + \alpha(y^{(i)} - \sum_{j=1}^n \beta_j \phi(x^{(j)})^T \phi(x^{(i)})), \forall i=1 \dots n \quad \downarrow \text{PSD}$$

$$\phi(x^{(i)})^T \phi(x^{(j)}) = \langle \phi(x^{(i)}), \phi(x^{(j)}) \rangle \triangleq K(x^{(i)}, x^{(j)}) = K(x, z) \in \mathbb{R}^{n \times n}$$

1) pre-compute  $\langle \phi(x^{(i)}), \phi(x^{(j)}) \rangle$  for all  $i, j$  pairs

2) can often be computed efficiently

$$\langle \phi(x), \phi(z) \rangle = 1 + \sum_{i=1}^d x_i z_i + \sum_{i,j=1}^d x_i x_j z_i z_j + \sum_{i,j,k=1}^d x_i x_j x_k z_i z_j z_k \\ = 1 + \langle x, z \rangle + \langle x, z \rangle^2 + \langle x, z \rangle^3$$

$$\beta := \beta + \alpha(\bar{y} - K\beta), \beta_i := \beta_i + \alpha(y^{(i)} - \sum_{j=1}^n \beta_j K(x^{(i)}, x^{(j)}))$$

$$\text{prediction: } \Theta^T \phi(x) = \sum_{i=1}^n \beta_i \phi(x^{(i)})^T \phi(x) = \sum_{i=1}^n \beta_i K(x^{(i)}, x)$$

$$K(x, z) = (x^T z)^2 = \sum_{i=1}^d x_i z_i \sum_{j=1}^d x_j z_j = \sum_i \sum_j (x_i x_j)(z_i z_j)$$

$$K(x, z) = (x^T z + c)^2 = \sum_{i,j=1}^d (x_i x_j)(z_i z_j) + \sum_{i=1}^d (\sqrt{z_i} x_i)(\sqrt{z_i} x_i) + c^2$$

$$\text{Gaussian } K(x, z) = \exp\left(-\frac{1}{2\sigma^2} \|x - z\|^2\right)$$

$$z^T K z = \sum_i \sum_j z_i z_j K_{ij} = \sum_i \sum_j z_i z_j \phi(x^{(i)})^T \phi(x^{(j)}) z_j \\ = \sum_i \sum_j z_i z_j \sum_k \phi_k(x^{(i)}) \phi_k(x^{(j)}) z_j \\ = \sum_k \sum_i \sum_j z_i z_j \phi_k(x^{(i)}) \phi_k(x^{(j)}) z_j \\ = \sum_k (\sum_i z_i \phi_k(x^{(i)}))^2 \geq 0$$

## Support Vector machines

$$y \in \{-1, 1\}, h_{w,b}(x) = g(w^T x + b), g(z) = 1 \text{ if } z \geq 0 \text{ and } -1 \text{ o.w.}$$

$$\text{functional margin } \hat{\gamma}^{(i)} = y^{(i)}(w^T x^{(i)} + b), \hat{\gamma} = \min_i \hat{\gamma}^{(i)}, \text{ normalize } \|w\| = 1$$

$$\text{geometric margin } \gamma^{(i)} = \frac{1}{\|w\|} (w^T x^{(i)} + b) y^{(i)}, \gamma = \min_i \gamma^{(i)}, w, b \text{ to max } \gamma$$

optimal margin classifier:

$$\#1 \max \gamma \text{ s.t. } y^{(i)}(w^T x^{(i)} + b) \geq \gamma, i=1 \dots n \text{ and } \|w\|=1$$

$$\#2 \max \hat{\gamma} / \|w\| \text{ s.t. } y^{(i)}(w^T x^{(i)} + b) \geq \hat{\gamma}, i=1 \dots n$$

$$\#3 \min \frac{1}{2} \|w\|^2 \text{ s.t. } y^{(i)}(w^T x^{(i)} + b) \geq 1, i=1 \dots n$$

$$\min \frac{1}{2} \|w\|^2 = \frac{1}{2} \left( \sum_{i=1}^n \alpha_i y^{(i)} x^{(i)} \right)^T \left( \sum_{j=1}^n \alpha_j y^{(j)} x^{(j)} \right) \quad \text{constraint:} \\ y^{(i)} \left( \sum_{j=1}^n \alpha_j y^{(j)} x^{(j)} \right)^T x^{(i)} \geq 1 \\ = \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y^{(i)} y^{(j)} x^{(i)}^T x^{(j)} \\ = \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y^{(i)} y^{(j)} \underbrace{x^{(i)}^T}_{\langle x^{(i)}, x^{(j)} \rangle} \underbrace{x^{(j)}}_{\langle x^{(i)}, x^{(j)} \rangle} \\ y^{(i)} \left( \sum_j \alpha_j y^{(j)} \langle x^{(i)}, x^{(j)} \rangle \right) \geq 1$$

## Neural Network

$$X = \begin{bmatrix} 1 & x^{(1)} & \dots & x^{(n)} \end{bmatrix}^T \in \mathbb{R}^d$$

$$Z = \begin{bmatrix} 1 & z^{[1]}(1) & \dots & z^{[1]}(n) \end{bmatrix}^T \in \mathbb{R}^{\# \text{nodes}}$$

$$\text{forward: } Z^{[1]} = W^{[1]} X + b^{[1]}$$

$$a^{[1]} = \sigma(Z^{[1]}) \quad \sigma(z) = \frac{1}{1+e^{-z}}$$

$$Z^{[2]} = W^{[2]} a^{[1]} + b^{[2]}$$

$$a^{[2]} = \sigma(Z^{[2]})$$

$$Z^{[3]} = W^{[3]} a^{[2]} + b^{[3]}$$

$$a^{[3]} = \sigma(Z^{[3]}) = \hat{y}$$

$$C_{1 \times 1} \quad C_{1 \times 1}$$

$$\mathcal{L}^{(i)} = -[y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log (1-\hat{y}^{(i)})]$$

$$J(\hat{y}, y) = \frac{1}{n} \sum_i \mathcal{L}^{(i)}$$

$$\text{backward: } \frac{\partial J}{\partial W^{[3]}} = \underbrace{\frac{\partial J}{\partial a^{[3]}}} \cdot \underbrace{\frac{\partial a^{[3]}}{\partial Z^{[3]}}} \cdot \underbrace{\frac{\partial Z^{[3]}}{\partial W^{[3]}}}$$

$$\frac{\partial J}{\partial W^{[2]}} = \underbrace{\frac{\partial J}{\partial Z^{[3]}}} \cdot \underbrace{\frac{\partial Z^{[3]}}{\partial a^{[2]}}} \cdot \underbrace{\frac{\partial a^{[2]}}{\partial Z^{[2]}}} \cdot \underbrace{\frac{\partial Z^{[2]}}{\partial W^{[2]}}}$$

$$\frac{\partial J}{\partial W^{[1]}} = \underbrace{\frac{\partial J}{\partial Z^{[2]}}} \cdot \underbrace{\frac{\partial Z^{[2]}}{\partial a^{[1]}}} \cdot \underbrace{\frac{\partial a^{[1]}}{\partial Z^{[1]}}} \cdot \underbrace{\frac{\partial Z^{[1]}}{\partial W^{[1]}}}$$

$$\frac{\partial \mathcal{L}^{(i)}}{\partial W^{[3]}} = - \left[ y^{(i)} \underbrace{\frac{\partial}{\partial W^{[3]}} \log(w^{[3]} a^{[2]} + b^{[3]})}_{a^{[3]}} + (1-y^{(i)}) \underbrace{\frac{\partial}{\partial W^{[3]}} \log(1-\sigma(w^{[3]} a^{[2]} + b^{[3]}))}_{1-a^{[3]}} \right]$$

$$= -(y^{(i)} - a^{[3]}) a^{[2]T}$$

$$\frac{\partial \mathcal{L}^{(i)}}{\partial W^{[2]}} = \underbrace{\frac{\partial \mathcal{L}^{(i)}}{\partial a^{[3]}}} \cdot \underbrace{\frac{\partial a^{[3]}}{\partial Z^{[3]}}} \cdot \underbrace{\frac{\partial Z^{[3]}}{\partial a^{[2]}}} \cdot \underbrace{\frac{\partial a^{[2]}}{\partial Z^{[2]}}} \cdot \underbrace{\frac{\partial Z^{[2]}}{\partial W^{[2]}}}$$

$$\frac{\partial \mathcal{L}^{(i)}}{\partial W^{[3]}} = \underbrace{\frac{\partial \mathcal{L}^{(i)}}{\partial Z^{[3]}}} \cdot \underbrace{\frac{\partial Z^{[3]}}{\partial W^{[3]}}} \quad \begin{matrix} W^{[3]T} \\ a^{[2]} \\ (1-a^{[2]}) \\ a^{[1]T} \end{matrix}$$

$$\frac{\partial \mathcal{L}^{(i)}}{\partial W^{[2]}} = W^{[3]T} * a^{[2]} (1-a^{[2]}) (a^{[2]} - y^{(i)}) a^{[1]T}$$

## K-means

given:  $x^{(1)} \dots x^{(n)} \in \mathbb{R}^d$ ,  $K = \# \text{clusters}$

1) randomly select  $\mu^{(1)}, \dots, \mu^{(K)} \in \mathbb{R}^d$

2) for  $i=1 \dots n$ :  
 $C^{(i)} = \arg \min_{j \in 1 \dots K} \|x^{(i)} - \mu_j\|^2$   $\leftarrow \text{assign}$

3) for  $j=1 \dots K$ :  
 $\mu^{(j)} = \frac{1}{|C_j|} \sum_{i \in C_j} x^{(i)}$   $\leftarrow \text{compute } \mu$

$$\mu^{(j)} = \frac{1}{|C_j|} \sum_{i \in C_j} x^{(i)}$$

$$C_j = \{i : C^{(i)} = j\}$$

$$J(C, \mu) = \sum_i \|x^{(i)} - \mu_{C^{(i)}}\|^2$$

$$\text{real-valued regression: } \mathcal{L}(y, \hat{y}) = \frac{1}{2} (\hat{y} - y)^2$$

$$\text{logistic regression: } \mathcal{L}(y, \hat{y}) = -(y \log \hat{y} + (1-y) \log (1-\hat{y}))$$

$$\text{softmax over } K: \mathcal{L}(y, \hat{y}) = -\sum_{j=1}^K \mathbb{1}_{\{y=j\}} \log \hat{y}_j$$

Jensen's inequality:

$f$  is convex and  $X$  is random,

$$\mathbb{E}[f(X)] \geq f(\mathbb{E}[X])$$

$$\text{Accuracy} = \frac{TP + TN}{\text{all}}$$

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{recall sensitivity} = \frac{TP}{TP + FN}$$

$$\text{Specificity} = \frac{TN}{TN + FP}$$

$$F1 = \frac{2TP}{2TP + FP + FN}$$

$$\text{true pos. rate} = \frac{TP}{TP + FN}$$

$$\text{false pos. rate} = \frac{FP}{TN + FP}$$

$$(1-x)(1-y) = 1 - x - y + xy$$

$$(1-x)(1-y) = 1 - x - y + xy$$

$$(1-x)(1-y) = 1 - x - y + xy$$

$$(1-x)(1-y) = 1 - x - y + xy$$

$$(1-x)(1-y) = 1 - x - y + xy$$

$$(1-x)(1-y) = 1 - x - y + xy$$

$$(1-x)(1-y) = 1 - x - y + xy$$

$$(1-x)(1-y) = 1 - x - y + xy$$

$$(1-x)(1-y) = 1 - x - y + xy$$

$$(1-x)(1-y) = 1 - x - y + xy$$

$$(1-x)(1-y) = 1 - x - y + xy$$

$$\binom{n}{k} = \frac{n!}{k!(n-k)!}$$

choose K

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial \mathbf{y}}{\partial x_1} \\ \frac{\partial \mathbf{y}}{\partial x_2} \\ \vdots \\ \frac{\partial \mathbf{y}}{\partial x_n} \end{bmatrix}$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial x} & \frac{\partial y_2}{\partial x} & \dots & \frac{\partial y_m}{\partial x} \end{bmatrix}$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_2}{\partial x_1} & \dots & \frac{\partial y_m}{\partial x_1} \\ \frac{\partial y_1}{\partial x_2} & \frac{\partial y_2}{\partial x_2} & \dots & \frac{\partial y_m}{\partial x_2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_1}{\partial x_n} & \frac{\partial y_2}{\partial x_n} & \dots & \frac{\partial y_m}{\partial x_n} \end{bmatrix}$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{X}} = \begin{bmatrix} \frac{\partial y}{\partial x_{11}} & \frac{\partial y}{\partial x_{12}} & \dots & \frac{\partial y}{\partial x_{1q}} \\ \frac{\partial y}{\partial x_{21}} & \frac{\partial y}{\partial x_{22}} & \dots & \frac{\partial y}{\partial x_{2q}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y}{\partial x_{p1}} & \frac{\partial y}{\partial x_{p2}} & \dots & \frac{\partial y}{\partial x_{pq}} \end{bmatrix}$$

maximum likelihood estimate  $L = \prod_{i=1}^n P(x_i|y_i)$

SVM (cont.)

Hinge loss  $L(z, y) = \max(0, 1 - zy)$

Lagrangian  $\mathcal{L}(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^n \alpha_i [y^{(i)}(w^T x^{(i)} + b) - 1]$

$\nabla_w \mathcal{L}(w, b, \alpha) = w - \sum_i \alpha_i y^{(i)} x^{(i)} = 0$

$\frac{\partial}{\partial b} \mathcal{L} = \sum_i \alpha_i y^{(i)} = 0$

$\mathcal{L}(w, b, \alpha) = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle$

$\max_w W(w) = \mathcal{L}(w, b, \alpha)$  s.t.  $\alpha_i \geq 0$ ,  $\sum_i \alpha_i y^{(i)} = 0$

solve for  $\alpha_i, b$

$$w^* = \sum_{i=1}^n \alpha_i y^{(i)} x^{(i)}$$

$$b^* = \frac{1}{2} \left( \max_{i: y^{(i)}=-1} w^T x^{(i)} + \min_{i: y^{(i)}=1} w^T x^{(i)} \right)$$

$$\text{prediction}_n = h_{w,b}(x) = g(w^T x + b)$$

$$= g(\sum_i \alpha_i y^{(i)} \langle x^{(i)}, x \rangle + b)$$

$$= g(\sum_i \alpha_i y^{(i)} \langle x^{(i)}, x \rangle + b)$$

regularization = penalty C-hyperparameter violations

$$\min_w \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i \text{ s.t. } y^{(i)}(w^T x^{(i)} + b) \geq 1 - \xi_i, \xi_i \geq 0$$

if  $y^{(i)}(w^T x^{(i)} + b) > 0 \Rightarrow$  classified correctly

$$\mathcal{L}(w, b, \xi, r) = \frac{1}{2} w^T w + C \sum_i \xi_i - \sum_i \alpha_i [y^{(i)}(w^T x^{(i)} + b) - 1 - \xi_i] - \sum_i r_i \xi_i$$

$$\max_w W(w) = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle$$

$$\text{s.t. } 0 \leq \alpha_i \leq C, \sum_i \alpha_i = 0, i=1 \dots n$$

$$\alpha_i = 0 \Rightarrow y^{(i)}(w^T x^{(i)} + b) \geq 1$$

$$\alpha_i = C \Rightarrow y^{(i)}(w^T x^{(i)} + b) \leq 1$$

$$0 < \alpha_i < C \Rightarrow y^{(i)}(w^T x^{(i)} + b) = 1$$

Condition	Expression	Numerator layout, i.e. by $\mathbf{y}$ and $\mathbf{x}^T$	Denominator layout, i.e. by $\mathbf{y}^T$ and $\mathbf{x}$
$a$ is not a function of $\mathbf{x}$	$\frac{\partial a}{\partial \mathbf{x}} =$	0	
$\mathbf{A}$ is not a function of $\mathbf{x}$	$\frac{\partial \mathbf{A}\mathbf{x}}{\partial \mathbf{x}} =$	$\mathbf{A}$	$\mathbf{A}^T$
$\mathbf{A}$ is not a function of $\mathbf{x}$	$\frac{\partial \mathbf{x}^T \mathbf{A}}{\partial \mathbf{x}} =$	$\mathbf{A}^T$	$\mathbf{A}$
$a$ is not a function of $\mathbf{x}$ , $\mathbf{u} = \mathbf{u}(\mathbf{x})$	$\frac{\partial a\mathbf{u}}{\partial \mathbf{x}} =$	$a \frac{\partial \mathbf{u}}{\partial \mathbf{x}}$	
$\mathbf{v} = \mathbf{v}(\mathbf{x}), \mathbf{u} = \mathbf{u}(\mathbf{x})$	$\frac{\partial \mathbf{v}\mathbf{u}}{\partial \mathbf{x}} =$	$\mathbf{v} \frac{\partial \mathbf{u}}{\partial \mathbf{x}} + \mathbf{u} \frac{\partial \mathbf{v}}{\partial \mathbf{x}}$	$\mathbf{v} \frac{\partial \mathbf{u}}{\partial \mathbf{x}} + \mathbf{u}^T \mathbf{v}$
$\mathbf{A}$ is not a function of $\mathbf{x}$ , $\mathbf{u} = \mathbf{u}(\mathbf{x})$	$\frac{\partial \mathbf{A}\mathbf{u}}{\partial \mathbf{x}} =$	$\mathbf{A} \frac{\partial \mathbf{u}}{\partial \mathbf{x}}$	$\frac{\partial \mathbf{u}}{\partial \mathbf{x}} \mathbf{A}^T$
$\mathbf{u} = \mathbf{u}(\mathbf{x}), \mathbf{v} = \mathbf{v}(\mathbf{x})$	$\frac{\partial (\mathbf{u} + \mathbf{v})}{\partial \mathbf{x}} =$	$\frac{\partial \mathbf{u}}{\partial \mathbf{x}} + \frac{\partial \mathbf{v}}{\partial \mathbf{x}}$	
$\mathbf{u} = \mathbf{u}(\mathbf{x})$	$\frac{\partial g(\mathbf{u})}{\partial \mathbf{x}} =$	$\frac{\partial g(\mathbf{u})}{\partial u} \frac{\partial u}{\partial \mathbf{x}}$	$\frac{\partial u}{\partial \mathbf{x}} \frac{\partial g(\mathbf{u})}{\partial u}$
$\mathbf{u} = \mathbf{u}(\mathbf{x})$	$\frac{\partial f(g(\mathbf{u}))}{\partial \mathbf{x}} =$	$\frac{\partial f(g)}{\partial g} \frac{\partial g(\mathbf{u})}{\partial u} \frac{\partial u}{\partial \mathbf{x}}$	$\frac{\partial u}{\partial \mathbf{x}} \frac{\partial g(\mathbf{u})}{\partial u} \frac{\partial f(g)}{\partial g}$

$$f(x) = \sum_{i=1}^n b_i x_i \rightarrow f(x) = b^T x$$

$$\frac{\partial f(x)}{\partial x_k} = \frac{\partial}{\partial x_k} \sum_{i=1}^n b_i x_i = b_k$$

$$\nabla_x b^T x = b$$

$$f(y) = x^T A x, A \in \mathbb{S}$$

$$= \sum_{i=1}^n \sum_{j=1}^n A_{ij} x_i x_j$$

$$\frac{\partial f(y)}{\partial x_k} = \frac{\partial}{\partial x_k} \sum_{i=1}^n \sum_{j=1}^n A_{ij} x_i x_j =$$

$$= \frac{\partial}{\partial x_k} \left[ \sum_{i \neq k} \sum_{j \neq k} A_{ij} x_i x_j + \sum_{i \neq k} A_{ik} x_i x_k + \sum_{j \neq k} A_{kj} x_k x_j + A_{kk} x_k^2 \right]$$

$$= \sum_{i \neq k} A_{ik} x_i + \sum_{j \neq k} A_{kj} x_j + 2A_{kk} x_k$$

$$= \sum_{i=1}^n A_{ik} x_i + \sum_{j=1}^n A_{kj} x_j$$

$$= 2 \sum_{i=1}^n A_{ki} x_i$$

$$\nabla_x x^T A x = 2 A x, A \in \mathbb{S}$$

Loss functions:

$$\text{Least square error: } \frac{1}{2} (y - z)^2$$

$$\text{Logistic: } \log(1 + \exp(-yz))$$

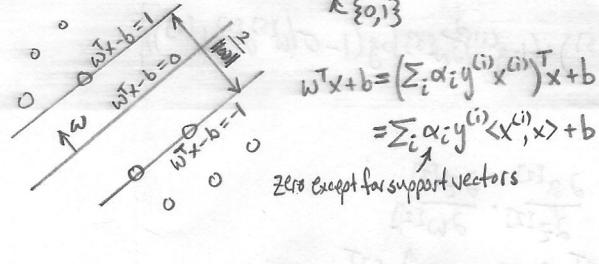
$$\text{Hinge: } \max(0, 1 - yz)$$

$$\text{Cross-entropy: } -[y \log(z) + (1-y) \log(1-z)]$$

$$J(\theta) = \sum_{i=1}^n L(h(x^{(i)}) y^{(i)})$$

$$\Theta \leftarrow \Theta - \alpha \nabla J(\Theta)$$

$$\text{Likelihood } L(\theta), \Theta^{\text{opt}} = \arg \max_{\theta} L(\theta)$$



regularization = penalty C-hyperparameter violations

$$\min_w \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i \text{ s.t. } y^{(i)}(w^T x^{(i)} + b) \geq 1 - \xi_i, \xi_i \geq 0$$

if  $y^{(i)}(w^T x^{(i)} + b) > 0 \Rightarrow$  classified correctly

$$\mathcal{L}(w, b, \xi, r) = \frac{1}{2} w^T w + C \sum_i \xi_i - \sum_i \alpha_i [y^{(i)}(w^T x^{(i)} + b) - 1 - \xi_i] - \sum_i r_i \xi_i$$

$$\max_w W(w) = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle$$

$$\text{s.t. } 0 \leq \alpha_i \leq C, \sum_i \alpha_i = 0, i=1 \dots n$$

$$\alpha_i = 0 \Rightarrow y^{(i)}(w^T x^{(i)} + b) \geq 1$$

$$\alpha_i = C \Rightarrow y^{(i)}(w^T x^{(i)} + b) \leq 1$$

$$0 < \alpha_i < C \Rightarrow y^{(i)}(w^T x^{(i)} + b) = 1$$