

## Classifications

- (1) Generative models (e.g., LDA) [We'll learn about LDA next lecture.]  
 - Assume sample points come from probability distributions, different for each class.  
 - Guess form of distributions  
 - For each class  $C$ , fit distribution parameters to class  $C$  points, giving  $f(X|C)$   
 - For each class  $C$ ,  $P(Y=C)$   
 - Bayes' Theorem gives  $P(Y|X)$   
 - If  $b=1$ , pick class  $C$  that maximizes  $P(Y=C|X=x)$   
 equivalently, maximizes  $\langle X | Y = C \rangle P(Y = C)$  [posterior probability]
- (2) Discriminative models (e.g., logistic regression)  
 [We'll learn about logistic regression in a few weeks.]  
 Model  $P(Y|X)$  directly
- (3) Decision boundary (e.g., SVM)  
 Model  $P(Y|X)$  directly (no posterior)
- Advantage of (1):  $P(Y|X)$  tells you probability; your guess is wrong  
 [This is sometimes true for SVMs, but not do.]
- Advantage of (1): you can discriminate earlier ( $X$ ) than it's small  
 Disadvantages of (1): often hard to estimate distributions accurately;  
 real distributions rarely match standard ones.

## Optimization Problems

- 1) Unconstrained:  $\min_w \max_u P(u)$   
 ↳ gradient descent
- 2) Constrained: Unconstraint + subject to  
 $g(w) = 0$   
 ↳ Lagrange
- 3) Linear Program:  $\min_w \max_u c^T w$  subject to  $Aw \leq b$   
 ↳ simplex
- 4) Quadratic Program:  $\min_w \max_u u^T Qw + c^T w$ ,  $Aw \leq b$   
 ↳ IP Q PSD,  $\rightarrow$  local min.

## Centroid

$$f(x) = (\mu_C - \mu_X) \cdot x - (\mu_C - \mu_X) \cdot \frac{\mu_C + \mu_X}{2}$$

## Perceptrons (+, -)

If  $x \cdot w \geq 0$ , predict 1

Else -1

$$\mathcal{L}(x, y) = \begin{cases} 0 & \text{if } y \cdot x \geq 0 \\ -y \cdot x \text{ otherwise} \end{cases} \quad * x \text{ is classifier's prediction.}$$

$$\hookrightarrow R(w) = \frac{1}{2} L(X; w, y)$$

$$= \frac{1}{2} \sum_i y_i x_i \cdot w$$

Find  $\min_w R(w)$  w/ Gradient Descent

$$\star O\left(\frac{\max\|X\|_2}{\gamma^2}\right) \text{ iterations}$$

where  $\gamma = \text{max margin}$

## GDA (QDA, LDA)

$$X \sim N(\mu, \sigma^2 I): f(x) = \frac{1}{(2\pi)^n \sigma^n} \exp\left(-\frac{\|x - \mu\|^2}{2\sigma^2}\right)$$

Predict  $\max_c P(X=x|Y=c) \pi_c$

↓ ln

$$\text{QDA: } \ln(\pi_1 \pi_2) f(x|\pi_1) = -\frac{\|x - \mu_1\|^2}{2\sigma_1^2} - \ln \pi_1 + \ln \pi_2$$

↪  $P(Y=c|x) = \frac{P(x|Y=c)\pi_c}{P(x|Y=c)\pi_c + P(x|Y=1)\pi_1}$

$$= S(Q_c \cdot \mathbf{1} \mathbf{1}^T - Q_1 \cdot \mathbf{1} \mathbf{1}^T)$$

LDA: Assume  $x \sim N$  with same variance  $S$ .

$$\text{Predict } \max_c \frac{\mu_c \cdot x}{S} - \frac{\|\mu_c\|^2}{2S} + \ln \pi_c$$

## Likelihood of Gaussian

$$L(\mu, \sigma^2; x_1, \dots, x_n) = f(x_1) \cdots f(x_n)$$

$$\hookrightarrow \ln L = \ln \prod_i f(x_i) = \frac{1}{2} \sum_i \|x_i - \mu\|^2$$

$$\rightarrow \hat{\mu} = \frac{1}{n} \sum_i x_i, \quad \hat{\sigma}^2 = \frac{1}{n} \sum_i \|x_i - \mu\|^2$$

For QDA:  $\forall C$ , find  $\hat{\mu}_c, \hat{\sigma}_c^2, \hat{\pi}_c: \frac{n_c}{n} \leq n_0$

$$\text{For LDA: } \hat{\sigma}^2 = \frac{1}{dn} \sum_c \sum_{i \in C} \|x_i - \mu\|^2$$

## SVMs

• Linear decision boundary not too close

### Hard Margin

$$\text{Constraints: } y_i (w \cdot x_i + b) \geq 1$$

$$\hookrightarrow \text{dist}(x_i, \text{hyperplane}) = \frac{|w|}{\sqrt{w^2}} \cdot |x_i| + \frac{|b|}{\sqrt{w^2}}$$

$$\therefore \text{margin is } \min_i \frac{1}{\sqrt{w^2}} |w \cdot x_i + b| \geq \frac{1}{\sqrt{w^2}} \rightarrow \text{slab width is } 2x$$

↪  $\min_w \|w\|^2 \text{ s.t. } y_i (w \cdot x_i + b) \geq 1$

### Soft Margin

• Allows some points to violate the margin

$$\hookrightarrow \text{Constraints: } y_i (w \cdot x_i + b) \geq 1 - \xi_i, \quad \xi_i \geq 0$$

$$\therefore \text{Find } \min_w \|w\|^2 + C \sum \xi_i$$

## Multivariate Gaussians

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

↓ PSD

$$\Sigma = V \Sigma^{\frac{1}{2}} V^T$$

$\Sigma^{\frac{1}{2}}$ :  $V \Sigma^{\frac{1}{2}} V^T$   $\star$  Maps sphere to ellipsoid

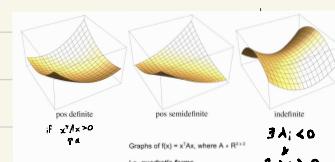
### QDA

$$\hat{z}_c = \frac{1}{n_c} \sum_{i \in C} (x_i - \hat{\mu}_c)(x_i - \hat{\mu}_c)^T$$

$$\text{Predict } \max_c -\frac{1}{2} (x - \hat{\mu}_c)^T \Sigma_c^{-1} (x - \hat{\mu}_c) - \frac{1}{2} \ln |\Sigma_c| + \ln \pi_c$$

### LDA

$$\text{Predict } \max_c \hat{\mu}_c \cdot x - \frac{1}{2} \hat{\mu}_c^T \hat{\mu}_c + \ln \pi_c$$



⑥ (p) Consider the two-dimensional bivariate normal distribution  $N(0, \Sigma)$  where the covariance matrix  $\Sigma$  is the matrix you learned in class and the mean is  $\mu = 0$ . Let  $f(x)$  be the PDF of that normal distribution, where  $x \in \mathbb{R}^2$ . What are the lengths of the major and minor axis of the ellipse?

$$f(x) = \frac{1}{2\pi\sqrt{|\Sigma|}} e^{-\frac{1}{2} x^T \Sigma^{-1} x}$$

Justify your answer.

The determinant of  $\Sigma$  is 10 (the product of the eigenvalues), or you can compute it the hard way). The normal PDF is

$$f(x) = \frac{1}{\sqrt{(2\pi)^2 |\Sigma|}} \exp\left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)\right) = \frac{1}{2\pi\sqrt{10}} \exp\left(-\frac{1}{2} x^T \Sigma^{-1} x\right)$$

Setting that to  $1/(4\pi\sqrt{10})$  gives

$$\begin{aligned} \frac{1}{2} x^T \Sigma^{-1} x &= -\ln \frac{1}{4} \\ \frac{1}{2} x^T \Sigma^{-1} x^2 &= -\ln 2 \\ |\Sigma^{-1}| x^2 &= -\ln 2 \end{aligned}$$

The eigenvalues of  $\Sigma^{-1}$  are  $1/\sqrt{5}$  and  $1/\sqrt{2}$ . The major axis of the ellipse is the eigenvector  $x$  that solves this equation and has eigenvalue  $1/\sqrt{5}$  (or eigenvalue 5 for  $\Sigma$ ); thus the major axis has length  $\sqrt{10 \ln 2}$ . The minor axis is the eigenvector  $x$  that solves this equation and has eigenvalue  $1/\sqrt{2}$  (or eigenvalue 2 for  $\Sigma$ ); thus the minor axis has length  $2\sqrt{\ln 2}$ .

## Bayesian Risk Minimization

$$\mathcal{R}(x) = E[L(c(x), V)]$$

$$= E[L(c(x), 1) P(V=1|X=x) + L(c(x), -1) P(V=-1|X=x)]$$

↓ If  $L$  symmetric, pick class that max posterior probability

## MATH Dangs

- Diagonalizable matrices  $\rightarrow V \Lambda V^T$
- If  $\Sigma$  is symmetric  $\rightarrow$  it's orthogonal ( $V^T V = I$ )
- For symmetric:  $A^T = V \Lambda V^T V^T = V \Lambda V^T$
- Only  $A^2$  is PSD
- $\text{Cor}(A, B) = E[(A - E[A])(B - E[B])^T]$
- $\text{Cor}(R, S) = E[(R - E[R])(S - E[S])^T] = E[R S^T] - E[R]E[S]^T$
- $(A - I\mathbf{v}) = 0$

• All  $\Sigma$ 's are PSD.

• for Gaussian dist

↪ ellipsoids radii of  $\sqrt{\lambda_i}$  of  $\Sigma$

• Ellipse in the form

• Gradient descent:  $w_{t+1} = w_t - \epsilon J'(w)$

## Regression

### Linear

$$\min_w \sum_{i=1}^n \|x_i^T w + b - y_i\|^2$$

write  $X$  as  $n \times (d+1)$  w/ bias

$$\rightarrow \min_w \|Xw - y\|^2 = RSS(w)$$

$$\|Xw - y\|^2 = w^T X^T X w - 2y^T X w + y^T y$$

$$VTRSS = 2X^T X w - 2X^T y = 0$$

$$\rightarrow X^T X w = X^T y$$

Solve: let  $w = (X^T X)^{-1} X^T y$   
pseudo-inverse

↳ if invertible, unique sol

### Weighted Least-Square Regression

Assign each  $x_i$ : a weight  $w_i$ , collect  $w_i$ 's in  $n \times n$  diag  $\Delta$

$$\min_w \|Xw - y\|^2 \Delta (Xw - y) = \sum_{i=1}^n w_i (X_i^T w - y_i)^2$$

$$\text{OR solve } X^T \Delta X w = X^T \Delta y.$$

### Newton's Method

At point  $v$ , approx  $J(w)$  near  $v$

$$\text{Taylor's: } \nabla J(w) = \nabla J(v) + (\nabla^2 J(v))(w-v) + O(\|w-v\|^2)$$

$$\rightarrow w = v - (\nabla^2 J(v))^{-1} \nabla J(v)$$

REPEAT:  $e \leftarrow \text{sol to } \nabla J^2(w), e = -\nabla J(w)$

$$w \leftarrow w + e$$

- doesn't know if max/min, expensive Hessian, ~~weights~~, non-smooth

+ If  $J$  quadratic  $\rightarrow$  1 step, right step length to reach min.

### Regularization

#### L2 (ridge)

$$F(w|x, y) = \frac{\ell(g(x, w) \cdot f(w))}{f(y|x)} = \frac{L(w) f(w)}{f(y|x)}$$

$$\rightarrow \min_w \|Xw - y\|^2 + \lambda \|w\|^2$$

↳ guarantees PSD, unique sol

$$\rightarrow \text{Normal eq: } (X^T X + \lambda I^T) w = X^T y$$

#### L1 (Lasso)

$$\rightarrow \min_w \|Xw - y\|^2 + \lambda \|w\|_1$$

Inclined to give  $w_i = 0$

#### Feature subset selection

All features increase var, but not all decreases bias.

#### Heuristic 1: forward step

↳ start w/ 0 features, keep adding best until err(validation)  $\uparrow$

↳  $O(d^2)$

#### Heuristic 2: backward step

↳ start w/ all, keep removing,  $O(d^2)$

### Logistic

$$\min_w \sum_{i=1}^n L(s(x_i \cdot w), y_i) \quad \# s = \frac{1}{1+e^{-y}} = \frac{g(x \cdot w)}{1-g(x \cdot w)}$$

$$= - \sum_{i=1}^n (y_i \ln(s(x \cdot w)) + (1-y_i) \ln(1-s(x \cdot w)))$$

$J(w)$  is convex  $\rightarrow$  gradient descent

$$\hookrightarrow \nabla J(w) = -X^T(y - s(Xw))$$

↳ Can use Newton's

$$\hookrightarrow \nabla^2 J(w) = \sum_{i=1}^n s_i(1-s_i) X_i X_i^T$$

$$= X^T \Delta X \text{ where } \Delta = \text{diag}(s_i(1-s_i))$$

$$\Rightarrow w \leftarrow w + (e: X^T \Delta X e = X^T(y - s))$$

### Statistical Justification

$$\text{Empirical risk: } \hat{R}(w) = \frac{1}{n} \sum_i L(h(x_i), y_i)$$

Suppose  $y_i \sim N(g(x_i), \sigma^2)$

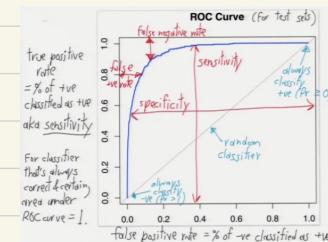
$$\hookrightarrow l(y_i; X, g) = -\frac{1}{2\sigma^2} (y_i - g(x_i))^2 - C$$

estimate  $g$  by Least Square

$$\cdot \hat{L}(h_i; X, g) = \prod_{i=1}^n h(x_i)^{y_i} (1-h(x_i))^{1-y_i}$$

$$\hookrightarrow \ln \hat{L}(h_i; X, g) = -\frac{1}{2} \sum_i \ln \text{logistic}(h(x_i), g)$$

Max likelihood  $\rightarrow$  minimize logistic loss

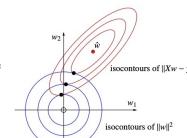
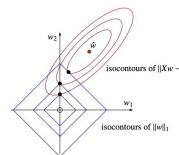


### Bias/Variance

$$\text{Bias}_0[\hat{f}(x; D)] = E_0[\hat{f}(x; D)] - f(x)$$

$$\text{Var}_0[\hat{f}(x; D)] = E_0[\hat{f}(x; D)^2] - E_0[\hat{f}(x; D)]^2$$

$$\begin{aligned} E[(y - \hat{f})^2] &= E[(f + e - \hat{f})^2] \\ &= E[(f - E[f])^2 + e^2 + 2(f - E[f])e] = E[(E[f] - \hat{f})^2] + 2E[(f - E[f])e] + 2E[(E[f] - f)(\hat{f} - E[\hat{f}])] \\ &= (f - E[f])^2 + E[e^2] = E[(E[f] - \hat{f})^2] + 2(f - E[f])E[e] + 2E[(f - E[f])(\hat{f} - E[\hat{f}])] \\ &= (f - E[f])^2 + E[e^2] = E[(E[f] - \hat{f})^2] \\ &= (f - E[f])^2 + \text{Var}[e] + \text{Var}[\hat{f}] \\ &= \text{Bias}[\hat{f}]^2 + \text{Var}[e] + \text{Var}[\hat{f}] \\ &= \text{Bias}[\hat{f}]^2 + \sigma^2 + \text{Var}[\hat{f}] \end{aligned}$$



### Support Vector Machine

start w/ all, keep removing,  $O(d^2)$