



SUPERVISED LEARNING

CLASSIFICATION

Keith Burghardt
USC Information Sciences Institute
DSCI 552 – Spring 2021

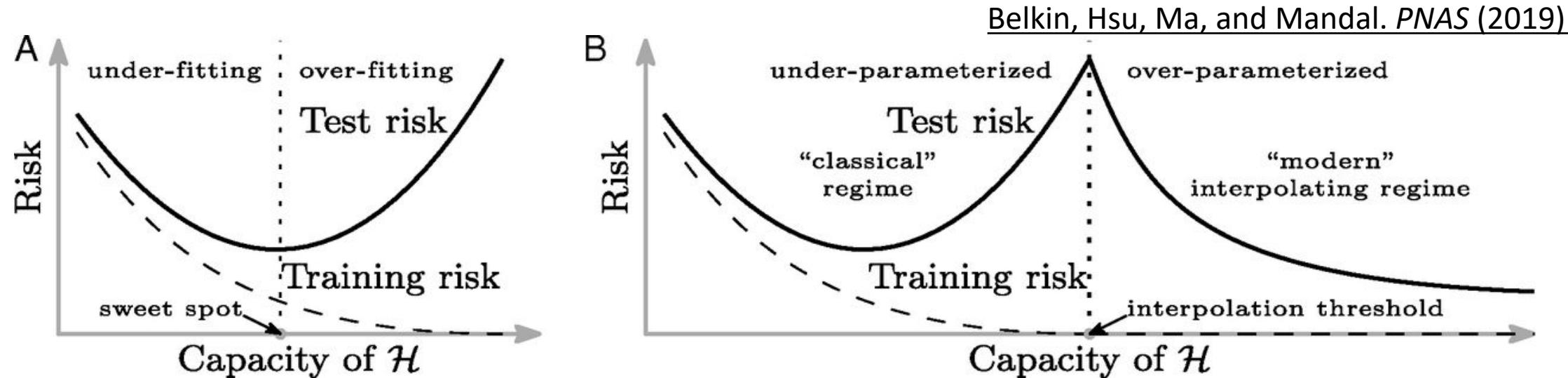


Topics this week

- **Reminder:**
 - Quiz 2 due Wednesday
 - Quiz 3 will be posted Wednesday
 - **Literature review and Quiz 3 are due February 10 at 10 a.m. PT (next Wednesday)**
- Supervised classification
 - Bayesian inference
 - Parametric models (ex: logistic regression)
 - K-Nearest Neighbo(u)rs
- Performance measures



Bias-Variance Curve Redux

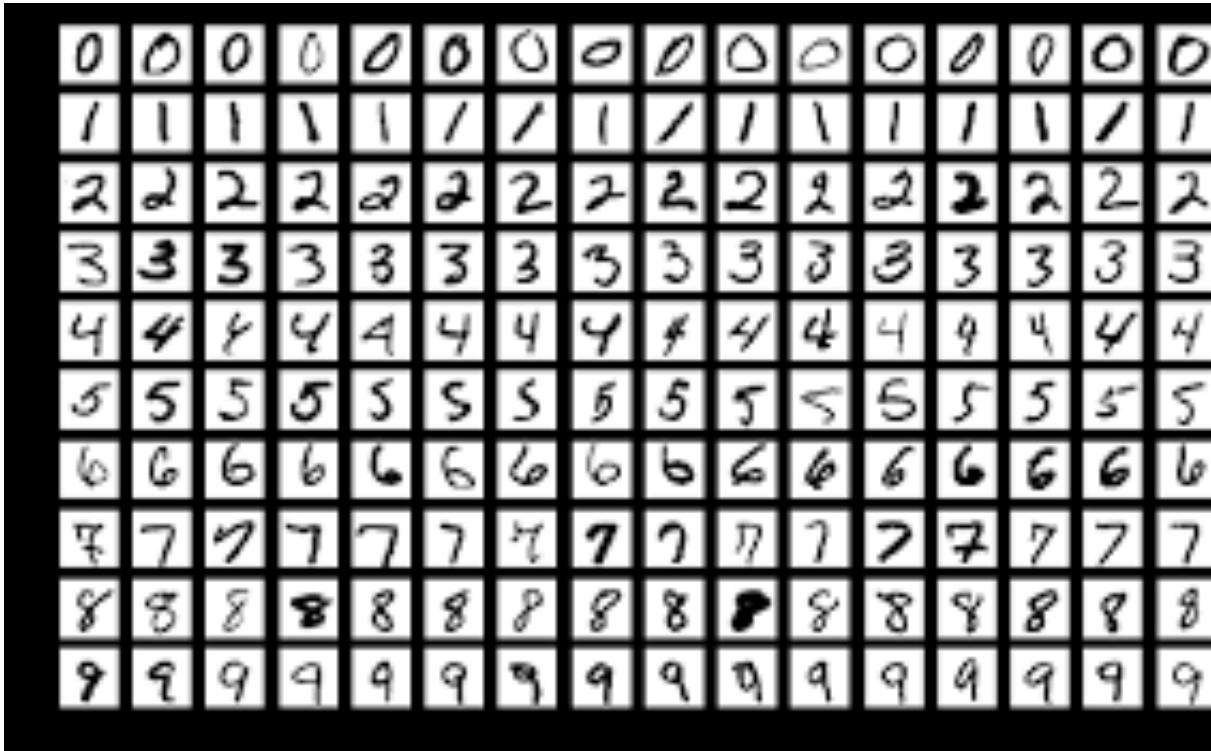


- Until 2019, the left-hand description was standard
- RQ: Why do neural networks *work*?
 - Old models = Dozens of parameters
 - NNs: thousands to TRILLIONS of parameters
- Empirically found that non-linear methods *improve* with parameters



Example: MNIST

- Train on 48,000 images; test on 28,000 images
- 32 x 32 pixels
- We can exactly fit every pixel with 43M parameters



VGG19

Accuracy: 99.5%

Parameters: 20M

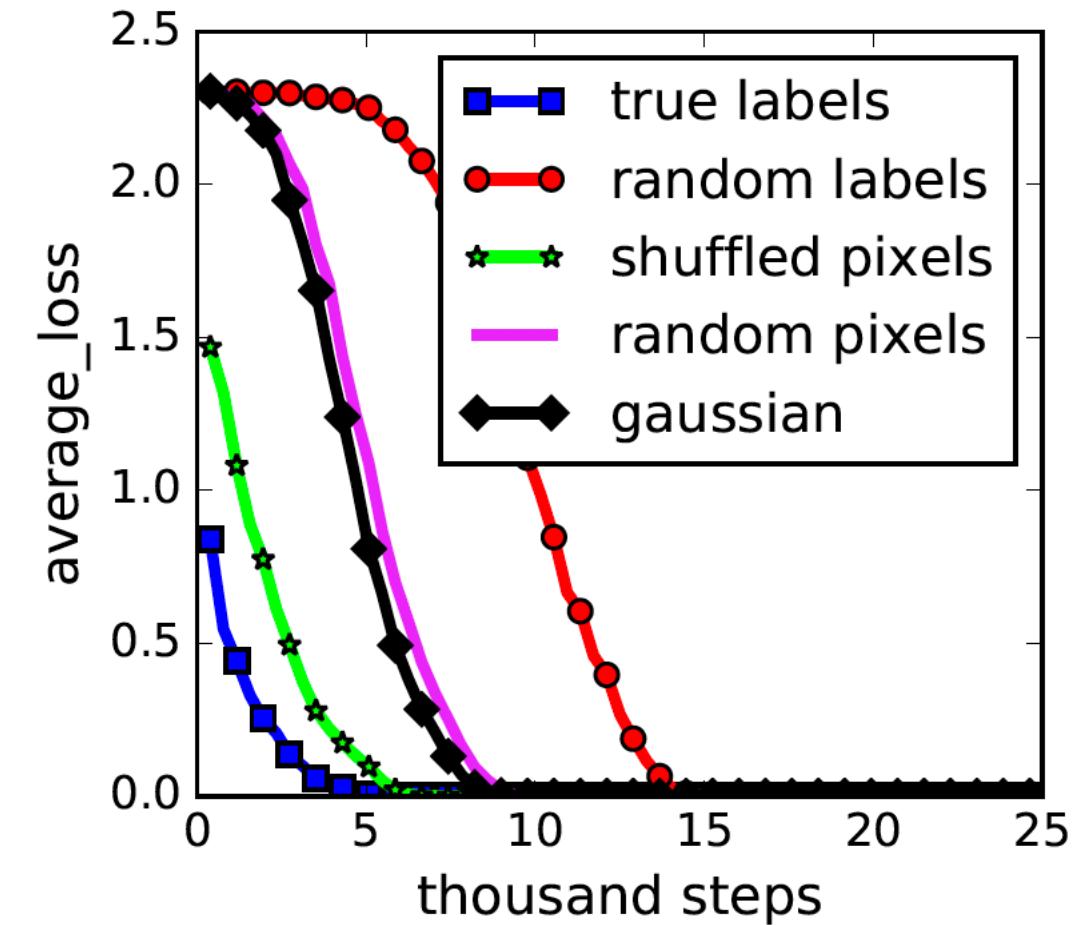
Only 2 pixels/parameter(!)

<https://www.kaggle.com/muerbingsha/mnist-vgg19>



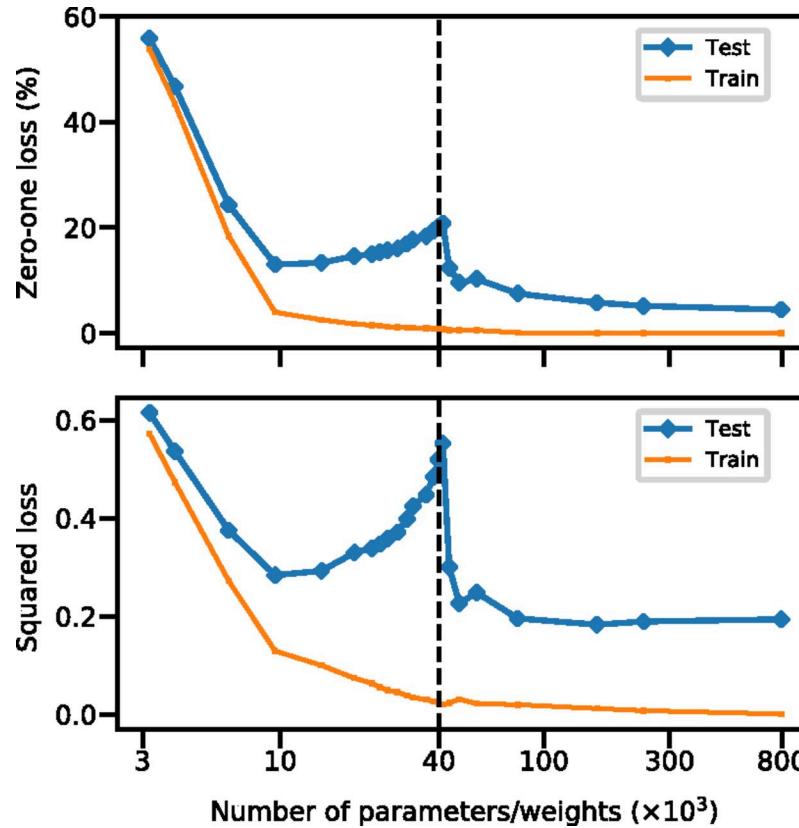
Training on Random Labels?

- Zhang et al., (2017) asked: what if we train on random labels?
- Neural networks can effectively memorize training data!
- They did not know why a neural network could “overfit” yet work so well

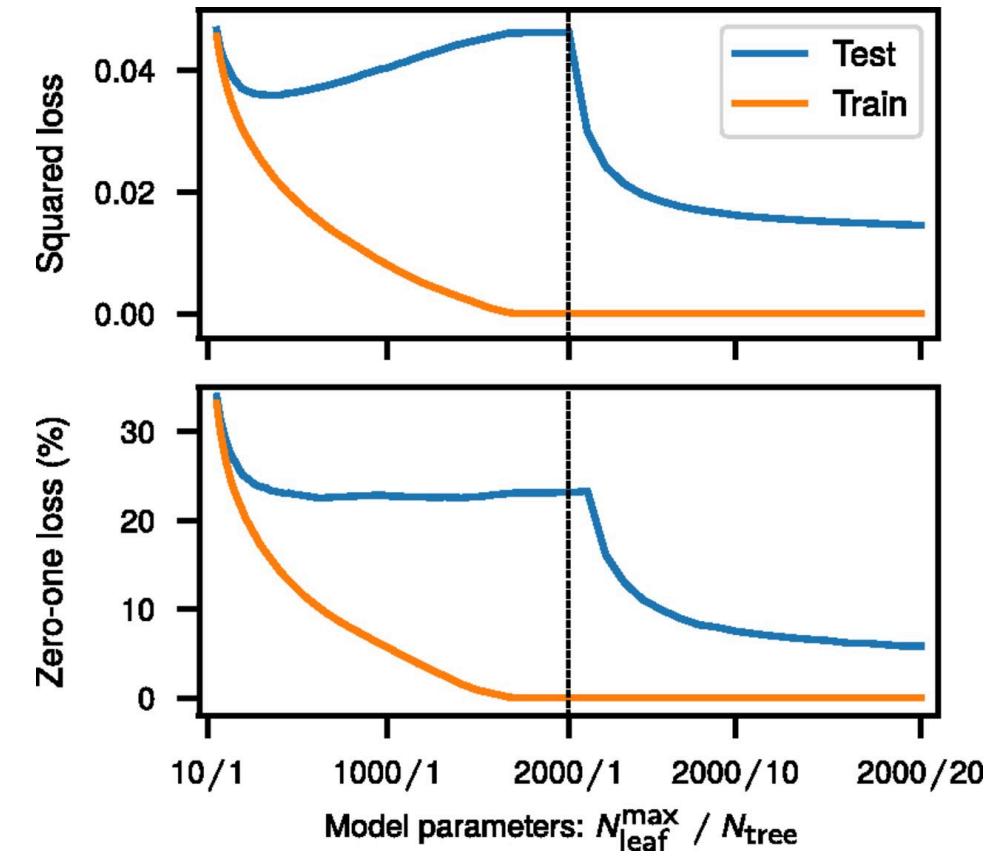




Examples on “MNIST”



Neural network



Random Forest



More Parameters = Better?

- Yes and no
- For many models/datasets, yes
- Imagine we have 3 datapoints, should we fit a 1000-parameter model on the data? No!
- Clearly, there is a limit to when this is correct
- Nonetheless, keep in mind: don't reduce the number of parameters just for the sake of it
- Use validation/testing to infer the best model for the job

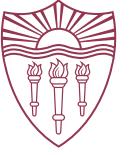


CLASSIFICATION MODELS



What we mean by classification

- The prior paper has real life examples, e.g., online check deposits:
 - Apps allow you to deposit a check by just taking a picture
 - Why can't you tell bank the check was for “\$1900” instead of “\$10.00”?
- Deposit checks at scale: banks need neural network to
 - distinguish “9” from “0”, and
 - “.” from smudge
- Classification: sort data into classes, e.g., 0,1,2,...
 - Supervised (learn what “0” means)
 - Unsupervised (find distinct clusters)
- Today, we will discuss how supervised classification works

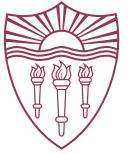


BAYESIAN INFERENCE



Let's simplify the problem: coin toss

2:40



“Call it”: What would you predict?

- Scenario:
 - You take a break from ML, and decide to run a gas station in Texas
 - Anton Chigurh arrives and asks you to guess the coin toss
 - You (correctly) suspect that losing the coin toss is not in your best interest
- What do you predict?
- If you never saw the coin before, you may assume it is “fair”
 - We call this the “*prior*”
- Given you saw the coin tossed a few times: {H,T,H,H,H,H,T}
 - We call this the “*evidence*”
- You calculate the *likelihood* of 2T and 5H given your estimated probability of heads, p



This is Bayesian Statistics!

- Let “p” be the probability of heads
 - This is a continuous variable
 - We are not looking at Bayesian inference with classes just yet
- What could “p” be?
- Let x be the evidence ($\#H, \#T$)
- Bayes rule:
 - $\Pr(p|x) = \Pr(p) \Pr(x|p)/\Pr(x)$
- The possible values of p depends on
 - Prior: $\Pr(p)$ (overall likelihood next flip is heads)
 - Likelihood: $\Pr(\#H,\#T|p)$ ($\#H,\#T$ if probability of heads is p)
 - Probability of evidence: $\Pr(\#H,\#T)$ ($=1$ in our case)

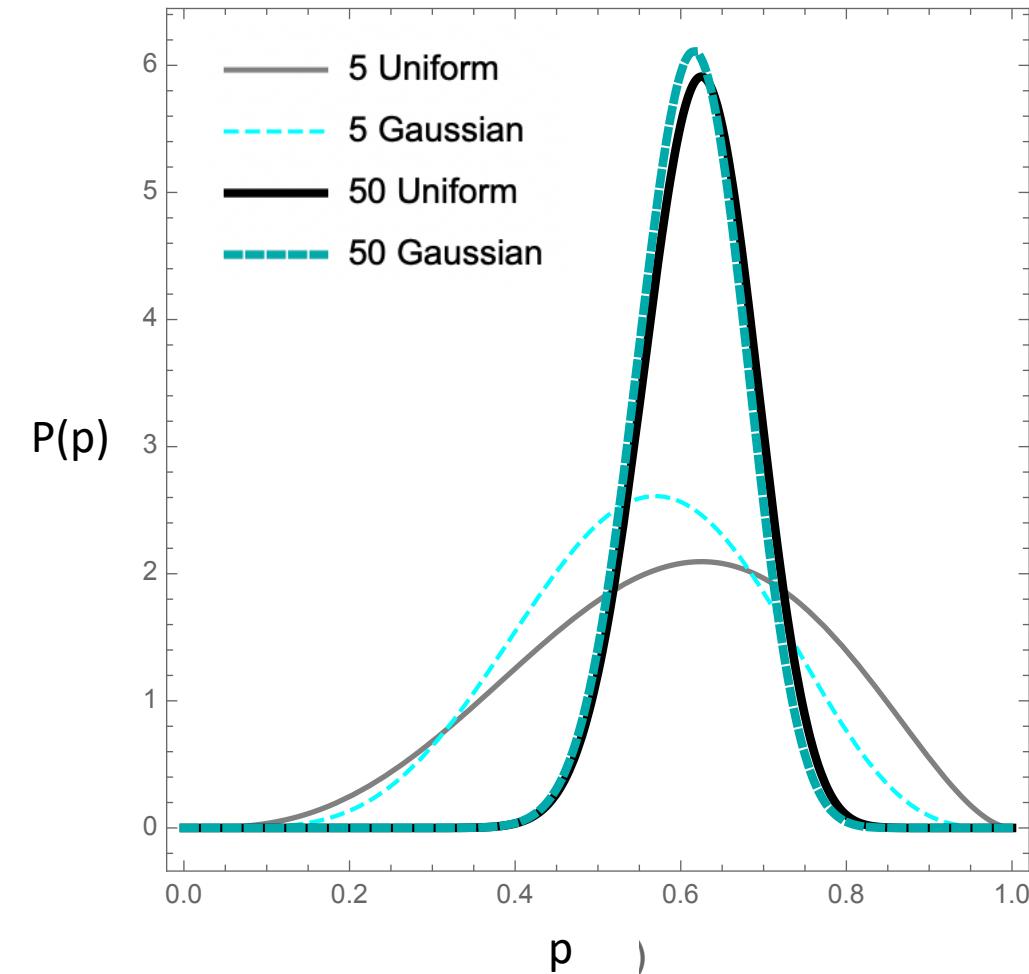
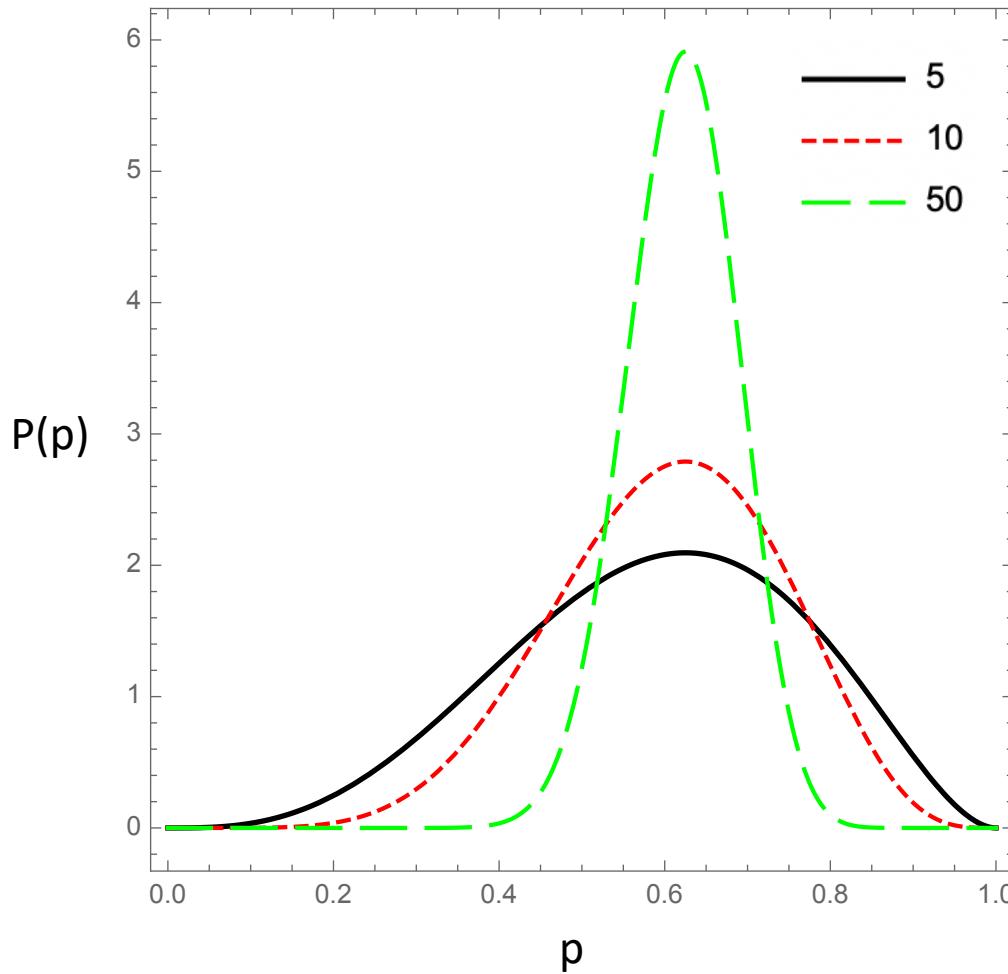


Bayes Rule

- How does Bayes' Rule come about?
- Conditional probability definition: $P(\theta|x) = P(\theta,x)/p(x)$
- Similarly, $P(x|\theta) = P(\theta,x)/p(\theta)$
- $P(\theta,x) = P(\theta)*P(x|\theta) = p(x)*P(\theta|x)$
- Therefore, divide by $p(x)$, and $P(\theta|x) = P(\theta,x)/p(x)$
- H/T example:
 - $P(x=H|p) = p; P(x=T|p) = 1-p$
 - $P(X|p) = p^x(1-p)^{1-x}$
 - Likelihood: $P(\{x_1,x_2,x_3,\dots\}|p) = p^{\#H} (1-p)^{\#T}$
 - Prior, e.g., uniform = $\frac{1}{2}$, or biased, e.g., $p^2 (1-p)^2$
 - If uniform prior: $P(p|x) \sim p^{\#H} (1-p)^{\#T}$



Effect of N, Prior





What we notice

- Maximum likelihood = most likely posterior if uniform prior
- Priors matter when data is small!
 - Posteriors show error reduces as more evidence is gathered
 - As more data is gathered, posteriors approximate true
- Useful in, e.g., neural networks to estimate classification error
- Real world example:
 - You or I could have really different views of the world
 - But if we agree on the facts, then as evidence is gathered we will approach the truth
 - E.g., Greek or Christian myths on how earth formed turned into our modern understanding of earth's formation as we gathered more evidence



“Call it.”

- What would you choose: heads or tails?
- We notice that evidence does not rule out a fair coin
- MLE: $p = 5/7$, so predict heads





Well done!

3:44



Bayes' Rule: K Classes

$$\begin{aligned} P(C_i | \mathbf{x}) &= \frac{p(\mathbf{x} | C_i)P(C_i)}{p(\mathbf{x})} \\ &= \frac{p(\mathbf{x} | C_i)P(C_i)}{\sum_{k=1}^K p(\mathbf{x} | C_k)P(C_k)} \end{aligned}$$

$$P(C_i) \geq 0 \text{ and } \sum_{i=1}^K P(C_i) = 1$$

choose C_i if $P(C_i | \mathbf{x}) = \max_k P(C_k | \mathbf{x})$



Bayes' Estimator for Classes

- Until now, I showed coin tossing
 - This is an example of estimating a continuous variable, p
- We can also estimate the probability of discrete classes
- Same rules as before:
 - Treat ϑ as a random var with prior $p(\vartheta)$
 - Bayes' rule: $p(\vartheta|\mathcal{X}) = p(\mathcal{X}|\vartheta) p(\vartheta) / p(\mathcal{X})$
 - Full: $p(x|\mathcal{X}) = \int p(x|\vartheta) p(\vartheta|\mathcal{X}) d\vartheta$
 - Maximum a Posteriori (MAP):
$$\vartheta_{\text{MAP}} = \operatorname{argmax}_{\vartheta} p(\vartheta|\mathcal{X})$$
 - Maximum Likelihood (ML): $\vartheta_{\text{ML}} = \operatorname{argmax}_{\vartheta} p(\mathcal{X}|\vartheta)$
 - Bayes': $\vartheta_{\text{Bayes'}} = E[\vartheta|\mathcal{X}] = \int \vartheta p(\vartheta|\mathcal{X}) d\vartheta$

Bayes' Estimator: Example

- $x^t \sim \mathcal{N}(\vartheta, \sigma_0^2)$ and $\vartheta \sim \mathcal{N}(\mu, \sigma^2)$

- $\vartheta_{\text{ML}} = m$

- $\vartheta_{\text{MAP}} = \vartheta_{\text{Bayes'}} =$

$$E[\vartheta | \mathcal{X}] = \frac{N/\sigma_0^2}{N/\sigma_0^2 + 1/\sigma^2} m + \frac{1/\sigma^2}{N/\sigma_0^2 + 1/\sigma^2} \mu$$

Note: this is $O(1/N)$ different than MLE



Parametric Classification

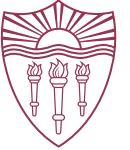
$$g_i(x) = p(x | C_i)P(C_i)$$

or

$$g_i(x) = \log p(x | C_i) + \log P(C_i)$$

$$p(x | C_i) = \frac{1}{\sqrt{2\pi}\sigma_i} \exp\left[-\frac{(x - \mu_i)^2}{2\sigma_i^2}\right]$$

$$g_i(x) = -\frac{1}{2} \log 2\pi - \log \sigma_i - \frac{(x - \mu_i)^2}{2\sigma_i^2} + \log P(C_i)$$



- Given the sample $\mathcal{X} = \{x^t, r^t\}_{t=1}^N$

$$x \in \Re \quad r_i^t = \begin{cases} 1 & \text{if } x^t \in C_i \\ 0 & \text{if } x^t \in C_j, j \neq i \end{cases}$$

- ML estimates are

$$\hat{P}(C_i) = \frac{\sum_t r_i^t}{N} \quad m_i = \frac{\sum_t x^t r_i^t}{\sum_t r_i^t} \quad s_i^2 = \frac{\sum_t (x^t - m_i)^2 r_i^t}{\sum_t r_i^t}$$

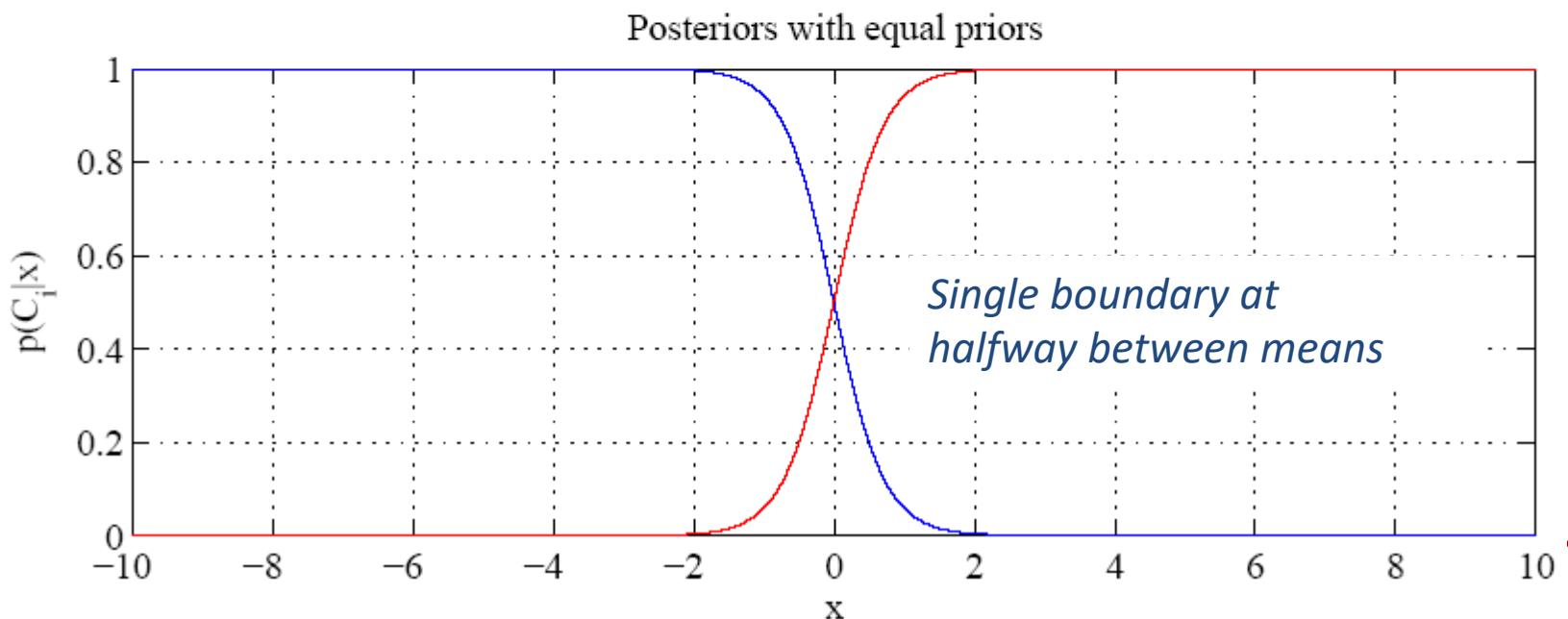
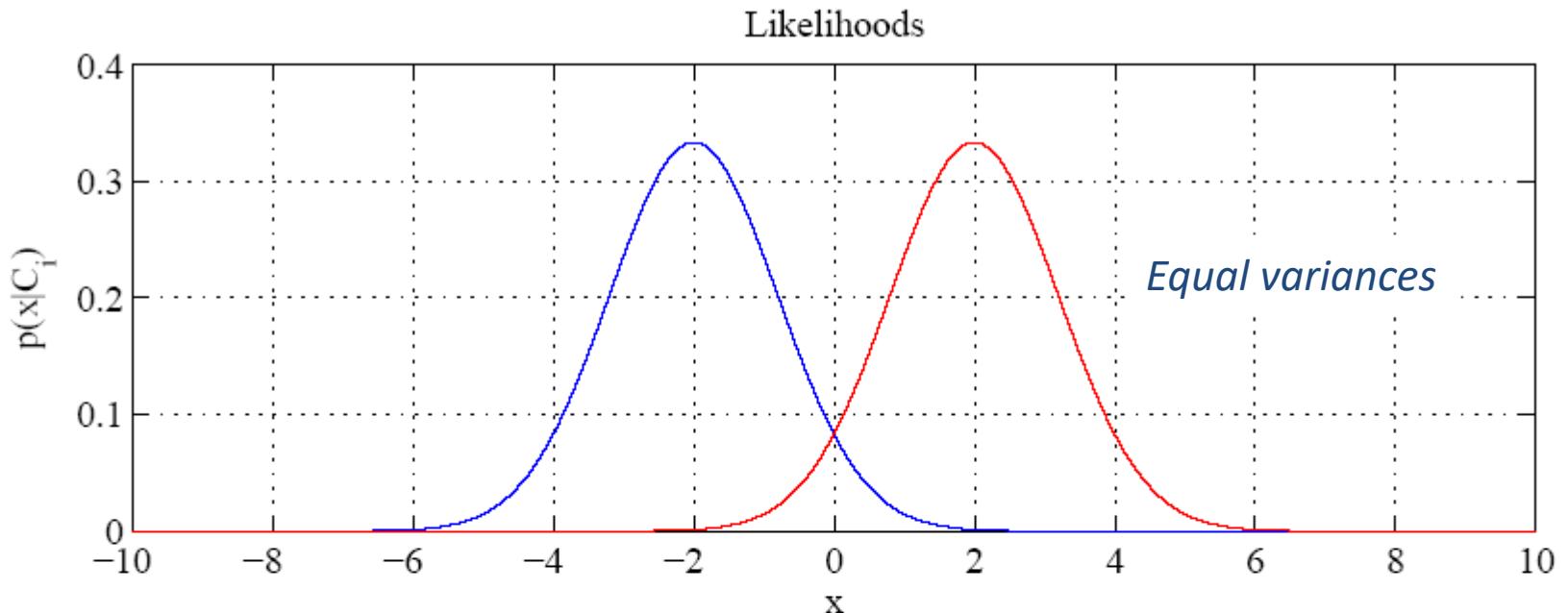
- Discriminant

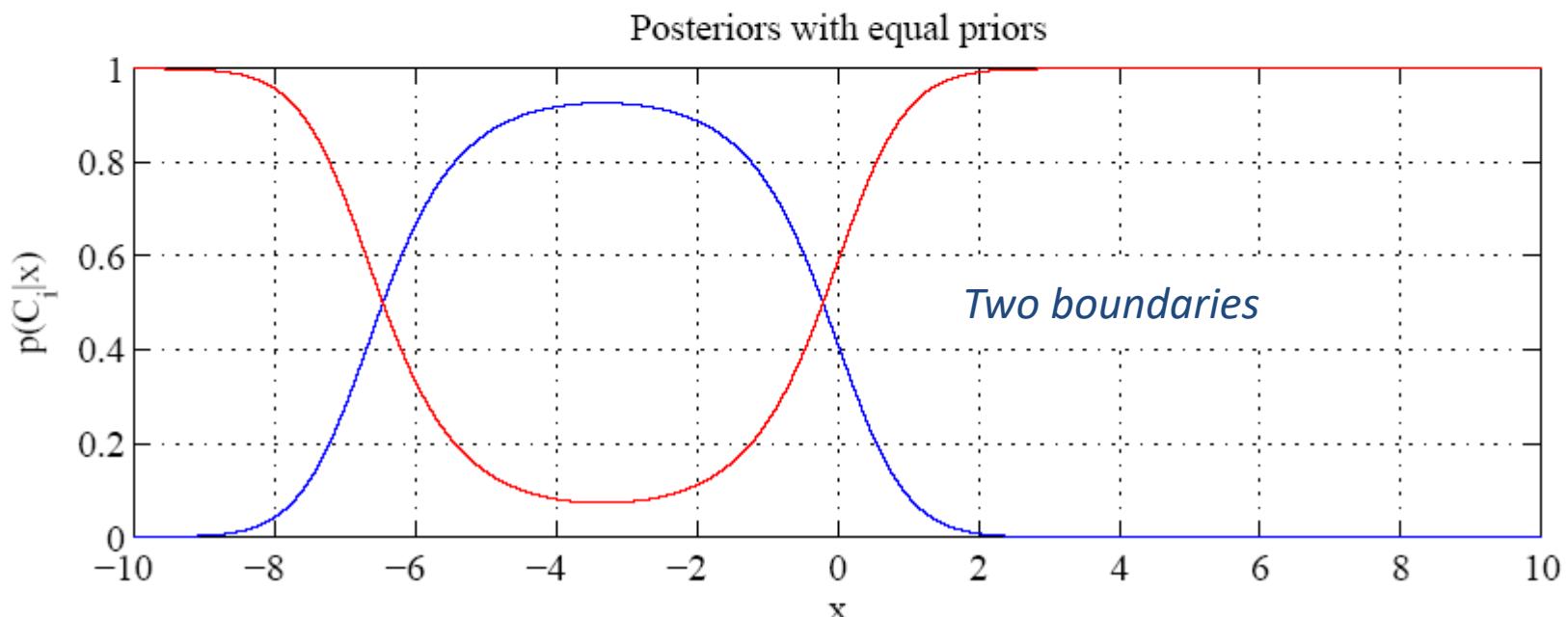
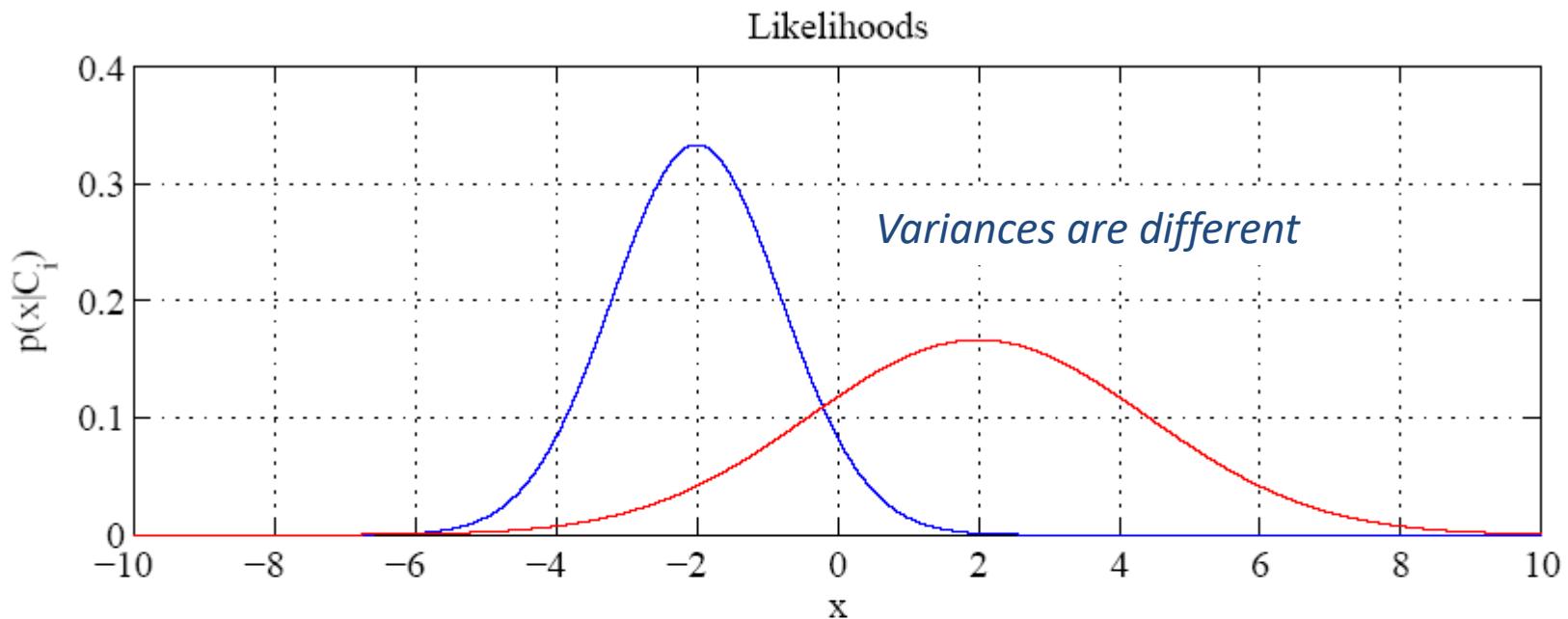
$$g_i(x) = -\frac{1}{2} \log 2\pi - \log s_i - \frac{(x - m_i)^2}{2s_i^2} + \log \hat{P}(C_i)$$



In English

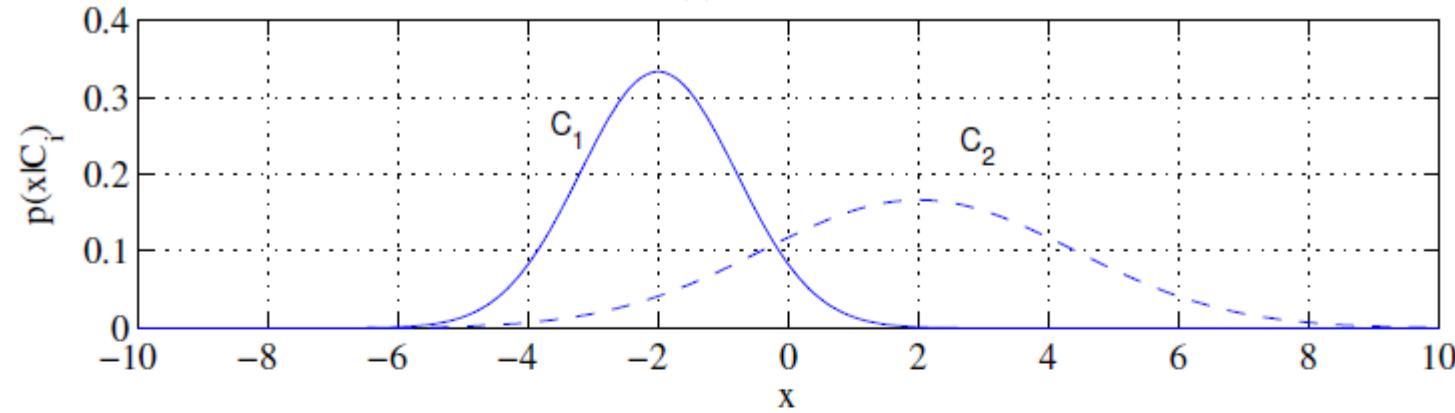
- Prior: how often each class appears
- Likelihood: how the data is distributed for each class
- Evidence: how often you see this set of feature values
- Posterior: Given you see a set of data, what is the likelihood of seeing class Ci given Ci's relative frequency (prior) and the frequency that the features would appear (likelihood)
- Discriminant is the unnormalized posterior
 - Whatever class has the largest discriminant has the largest posterior probability



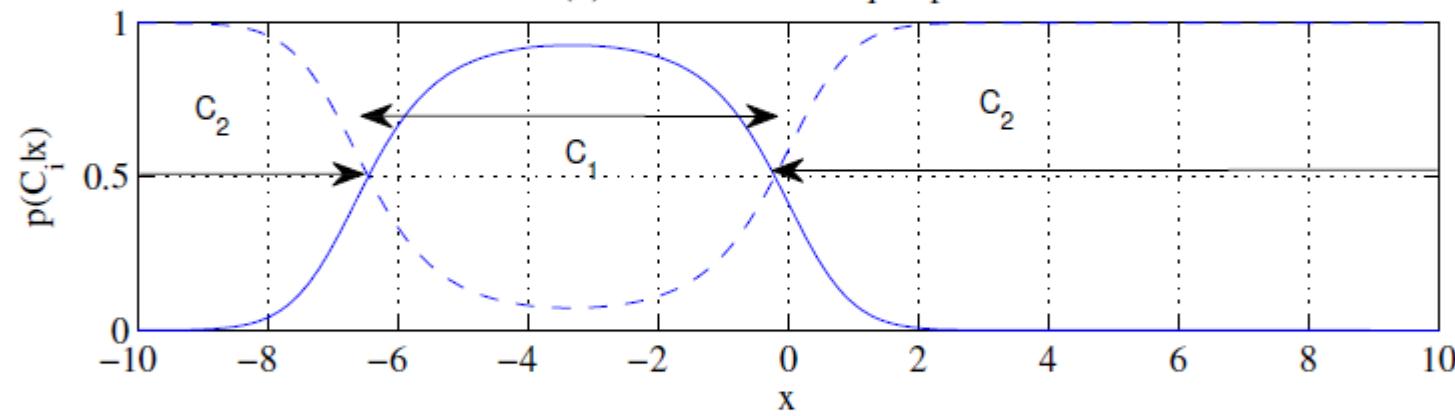




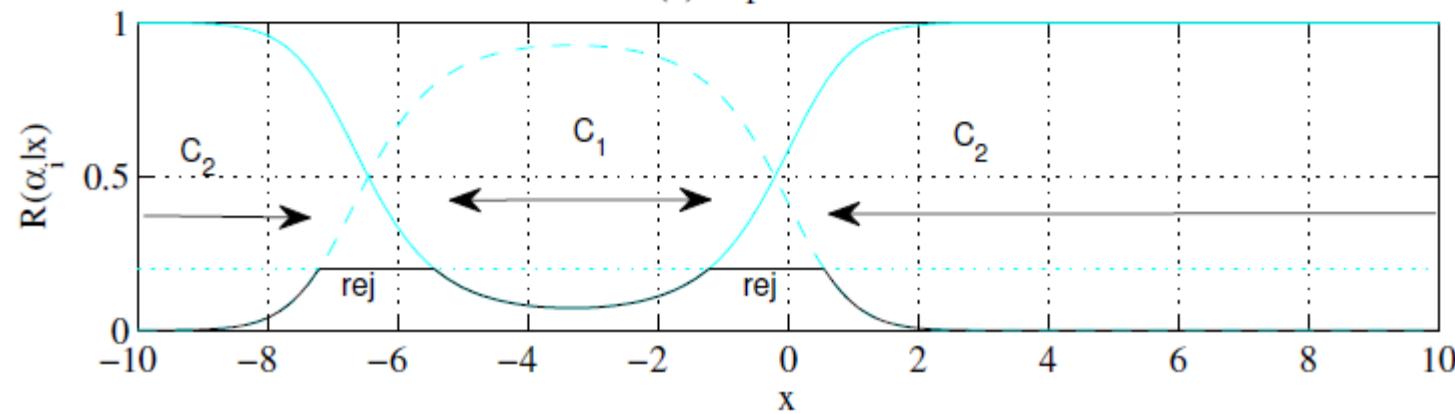
(a) Likelihoods



(b) Posteriors with equal priors



(c) Expected risks





MULTIVARIATE DATA



Multivariate Data

- Multiple measurements (sensors)
- d inputs/features/attributes: d -variate
- N instances/observations/examples

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^1 & \mathbf{x}_2^1 & \dots & \mathbf{x}_d^1 \\ \mathbf{x}_1^2 & \mathbf{x}_2^2 & \dots & \mathbf{x}_d^2 \\ \vdots & & & \\ \mathbf{x}_1^N & \mathbf{x}_2^N & \dots & \mathbf{x}_d^N \end{bmatrix}$$



Multivariate Parameters

Mean: $E[\mathbf{x}] = \boldsymbol{\mu} = [\mu_1, \dots, \mu_d]^T$

Covariance: $\sigma_{ij} \equiv \text{Cov}(x_i, x_j)$

Correlation: $\text{Corr}(x_i, x_j) \equiv \rho_{ij} = \frac{\sigma_{ij}}{\sigma_i \sigma_j}$

$$\Sigma \equiv \text{Cov}(\mathbf{X}) = E[(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^T] = \begin{bmatrix} \sigma_1^2 & \sigma_{12} & \cdots & \sigma_{1d} \\ \sigma_{21} & \sigma_2^2 & \cdots & \sigma_{2d} \\ \vdots & & & \\ \sigma_{d1} & \sigma_{d2} & \cdots & \sigma_d^2 \end{bmatrix}$$

Parameter Estimation

Sample mean \mathbf{m} : $m_i = \frac{\sum_{t=1}^N x_i^t}{N}, i = 1, \dots, d$

Covariance matrix \mathbf{S} : $s_{ij} = \frac{\sum_{t=1}^N (x_i^t - m_i)(x_j^t - m_j)}{N}$

Correlation matrix \mathbf{R} : $r_{ij} = \frac{s_{ij}}{s_i s_j}$

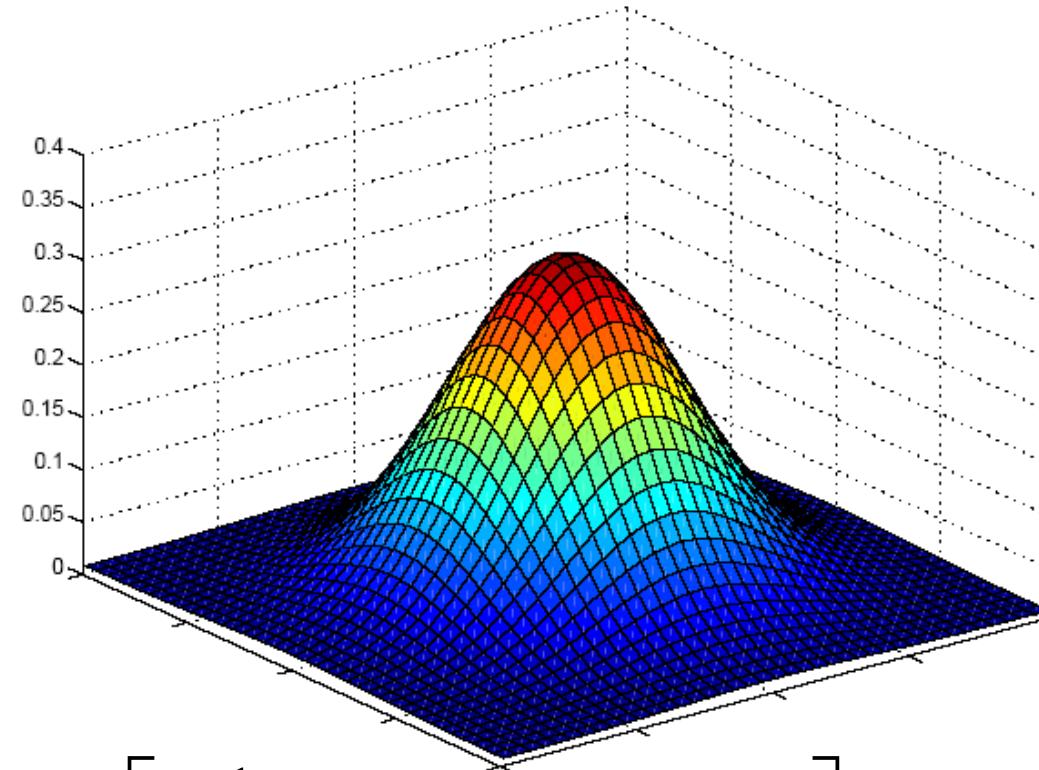


Estimation of Missing Values

- What to do if certain instances have missing attributes?
- Ignore those instances: not a good idea if the sample is small
- Use ‘missing’ as an attribute: may give information
- Imputation: Fill in the missing value
 - Mean imputation: Use the most likely value (e.g., mean)
 - Imputation by regression: Predict based on other attributes



Multivariate Normal Distribution



$$\mathbf{x} \sim \mathcal{N}_d(\boldsymbol{\mu}, \Sigma)$$

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



Multivariate Normal Distribution

- Mahalanobis distance: $(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})$
measures the distance from \mathbf{x} to $\boldsymbol{\mu}$ in terms of $\boldsymbol{\Sigma}$
(normalizes for difference in variances and correlations)
- Bivariate: $d = 2$

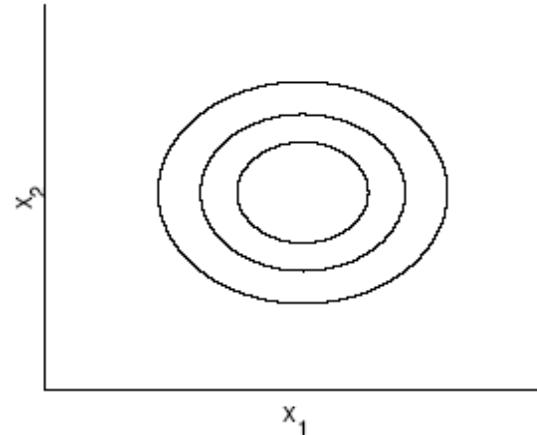
$$\boldsymbol{\Sigma} = \begin{bmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix}$$

$$p(x_1, x_2) = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \exp\left[-\frac{1}{2(1-\rho^2)}(z_1^2 - 2\rho z_1 z_2 + z_2^2)\right]$$
$$z_i = (x_i - \mu_i)/\sigma_i$$

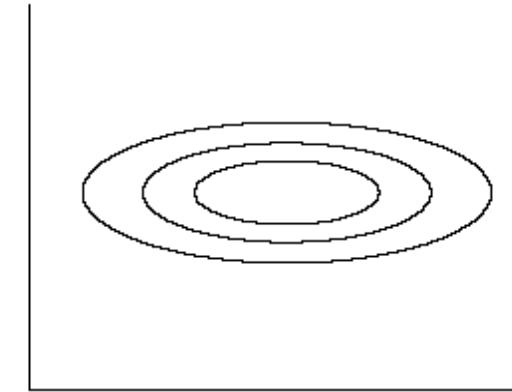


Bivariate Normal

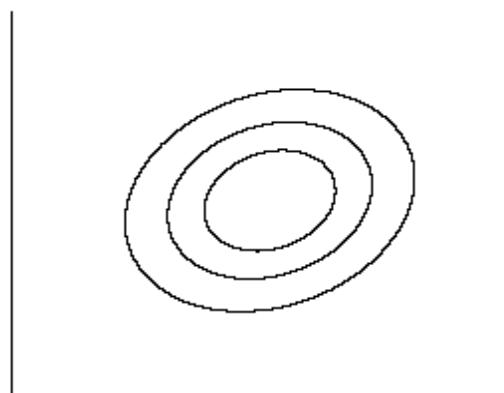
$\text{Cov}(x_1, x_2) = 0, \text{Var}(x_1) = \text{Var}(x_2)$



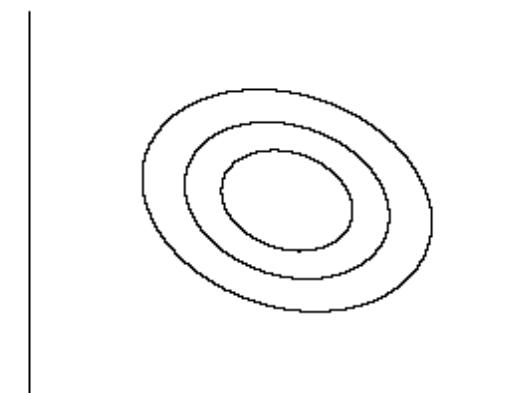
$\text{Cov}(x_1, x_2) = 0, \text{Var}(x_1) > \text{Var}(x_2)$

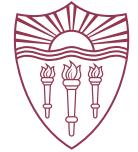


$\text{Cov}(x_1, x_2) > 0$



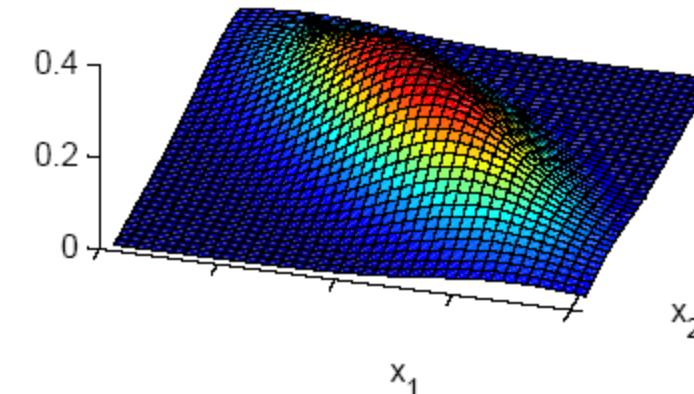
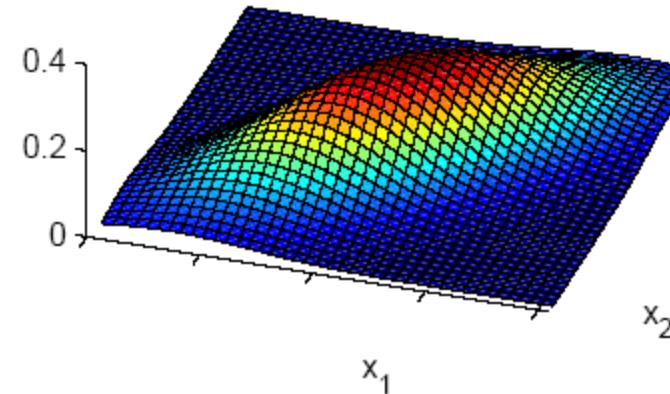
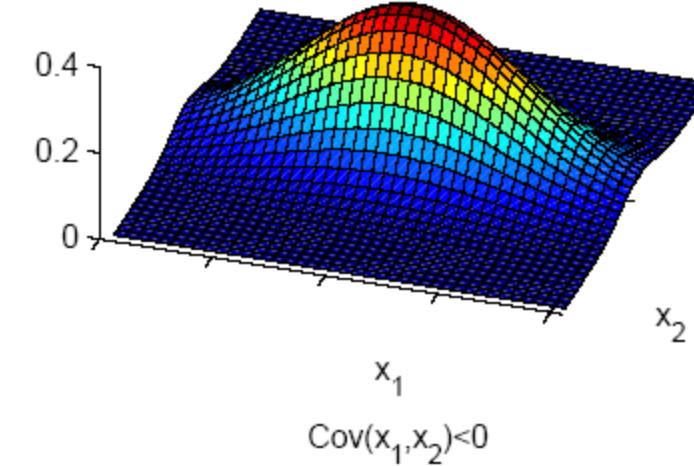
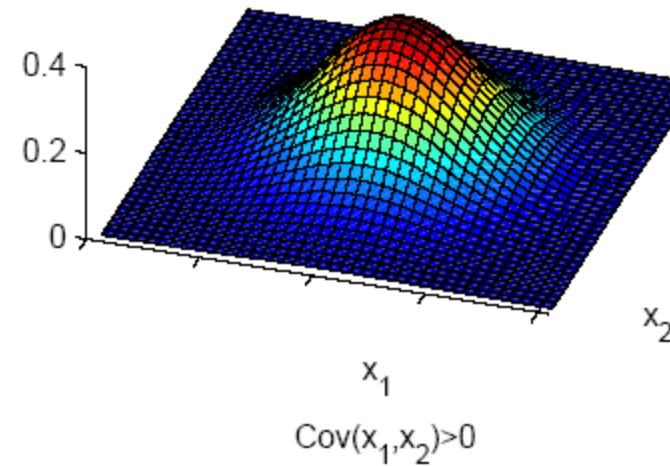
$\text{Cov}(x_1, x_2) < 0$





$\text{Cov}(x_1, x_2) = 0, \text{Var}(x_1) = \text{Var}(x_2)$

$\text{Cov}(x_1, x_2) = 0, \text{Var}(x_1) > \text{Var}(x_2)$





Independent Inputs: Naive Bayes

- If x_i are independent, offdiagonals of Σ are 0, Mahalanobis distance reduces to weighted (by $1/\sigma_i$) Euclidean distance:

$$p(\mathbf{x}) = \prod_{i=1}^d p_i(x_i) = \frac{1}{(2\pi)^{d/2} \prod_{i=1}^d \sigma_i} \exp \left[-\frac{1}{2} \sum_{i=1}^d \left(\frac{x_i - \mu_i}{\sigma_i} \right)^2 \right]$$

- If variances are also equal, reduces to Euclidean distance

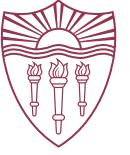
Multivariate Parametric Classification

- If $p(\mathbf{x} | C_i) \sim N(\boldsymbol{\mu}_i, \Sigma_i)$

$$p(\mathbf{x} | C_i) = \frac{1}{(2\pi)^{d/2} |\Sigma_i|^{1/2}} \exp\left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^T \Sigma_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i)\right]$$

- Discriminant functions

$$\begin{aligned} g_i(\mathbf{x}) &= \log p(\mathbf{x} | C_i) + \log P(C_i) \\ &= -\frac{d}{2} \log 2\pi - \frac{1}{2} \log |\Sigma_i| - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^T \Sigma_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) + \log P(C_i) \end{aligned}$$



Estimation of Parameters

$$\hat{P}(C_i) = \frac{\sum_t r_i^t}{N}$$

$$\mathbf{m}_i = \frac{\sum_t r_i^t \mathbf{x}^t}{\sum_t r_i^t}$$

$$\mathbf{S}_i = \frac{\sum_t r_i^t (\mathbf{x}^t - \mathbf{m}_i)(\mathbf{x}^t - \mathbf{m}_i)^T}{\sum_t r_i^t}$$

$$g_i(\mathbf{x}) = -\frac{1}{2} \log |\mathbf{S}_i| - \frac{1}{2} (\mathbf{x} - \mathbf{m}_i)^T \mathbf{S}_i^{-1} (\mathbf{x} - \mathbf{m}_i) + \log \hat{P}(C_i)$$

Different \mathbf{S}_i

- Quadratic discriminant

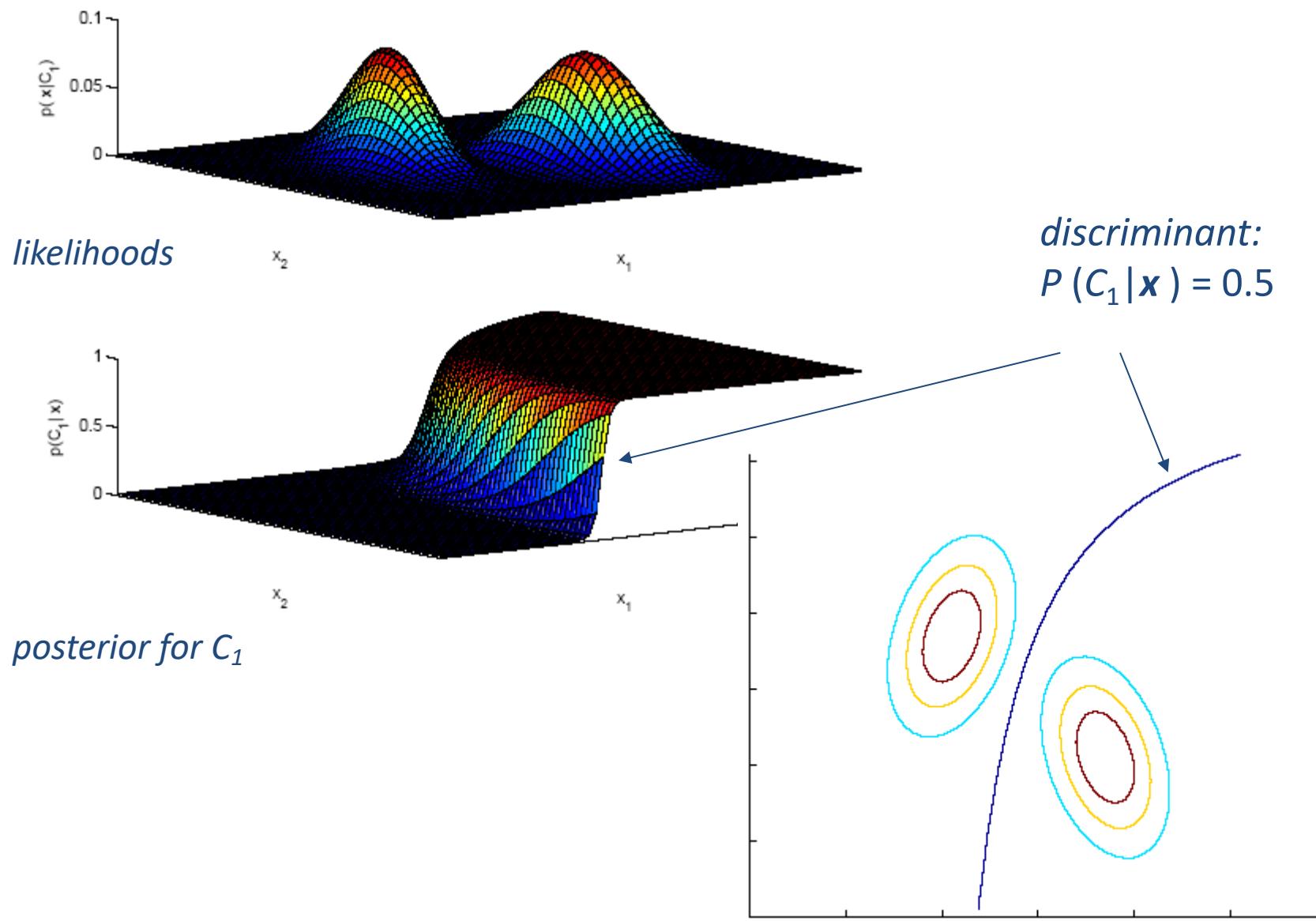
$$\begin{aligned}g_i(\mathbf{x}) &= -\frac{1}{2} \log |\mathbf{S}_i| - \frac{1}{2} \left(\mathbf{x}^T \mathbf{S}_i^{-1} \mathbf{x} - 2 \mathbf{x}^T \mathbf{S}_i^{-1} \mathbf{m}_i + \mathbf{m}_i^T \mathbf{S}_i^{-1} \mathbf{m}_i \right) + \log \hat{P}(C_i) \\&= \mathbf{x}^T \mathbf{W}_i \mathbf{x} + \mathbf{w}_i^T \mathbf{x} + w_{i0}\end{aligned}$$

where

$$\mathbf{W}_i = -\frac{1}{2} \mathbf{S}_i^{-1}$$

$$\mathbf{w}_i = \mathbf{S}_i^{-1} \mathbf{m}_i$$

$$w_{i0} = -\frac{1}{2} \mathbf{m}_i^T \mathbf{S}_i^{-1} \mathbf{m}_i - \frac{1}{2} \log |\mathbf{S}_i| + \log \hat{P}(C_i)$$





Common Covariance Matrix \mathbf{S}

- Shared common sample covariance \mathbf{S}

$$\mathbf{S} = \sum_i \hat{P}(C_i) \mathbf{S}_i$$

- Discriminant reduces to

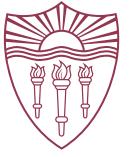
$$g_i(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \mathbf{m}_i)^T \mathbf{S}^{-1}(\mathbf{x} - \mathbf{m}_i) + \log \hat{P}(C_i)$$

which is a linear discriminant

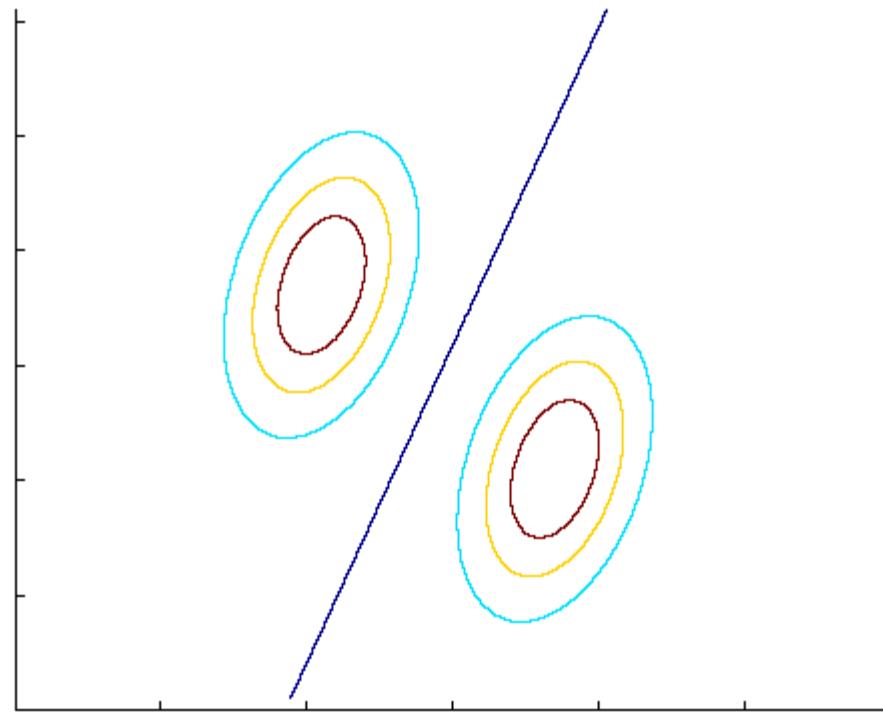
$$g_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + w_{i0}$$

where

$$\mathbf{w}_i = \mathbf{S}^{-1} \mathbf{m}_i \quad w_{i0} = -\frac{1}{2} \mathbf{m}_i^T \mathbf{S}^{-1} \mathbf{m}_i + \log \hat{P}(C_i)$$



Common Covariance Matrix \mathbf{S}





Diagonal \mathbf{S}

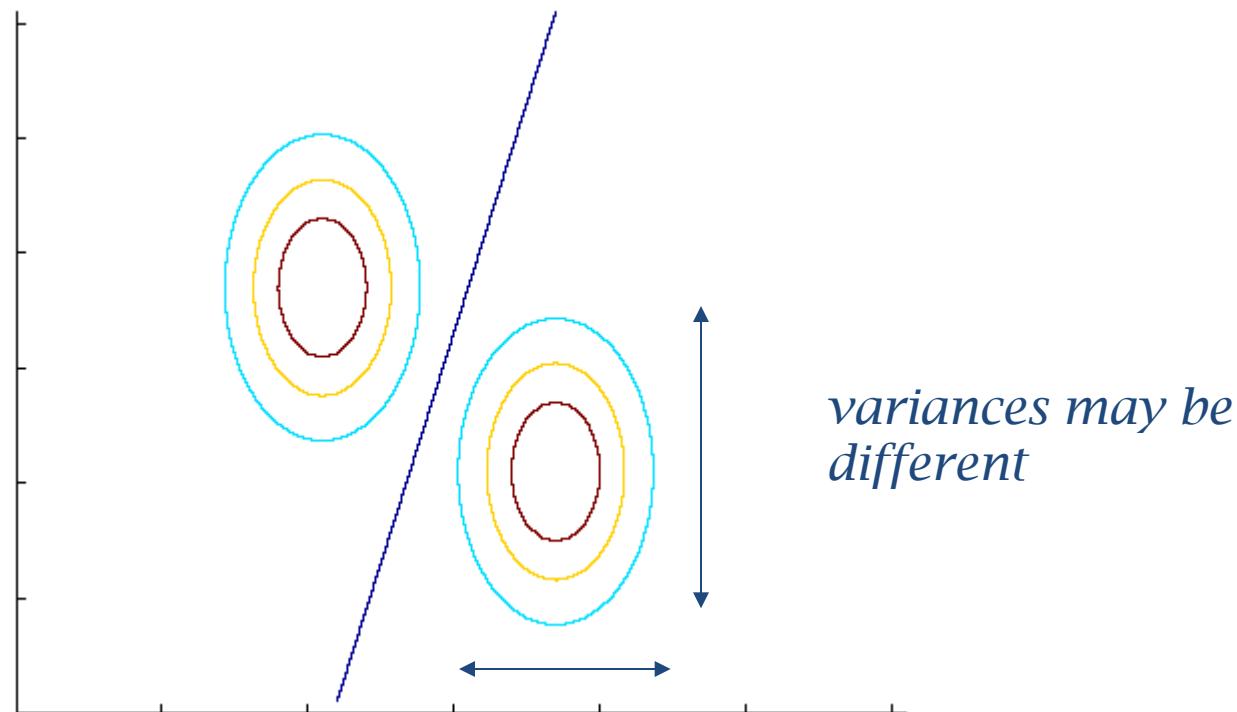
- When $x_j j = 1,..d$, are independent, Σ is diagonal
 $p(\mathbf{x}|C_i) = \prod_j p(x_j | C_i)$ (Naive Bayes' assumption)

$$g_i(\mathbf{x}) = -\frac{1}{2} \sum_{j=1}^d \left(\frac{x_j^t - m_{ij}}{s_j} \right)^2 + \log \hat{P}(C_i)$$

Classify based on weighted Euclidean distance (in s_j units) to the nearest mean



Diagonal S





Diagonal \mathbf{S} , equal variances

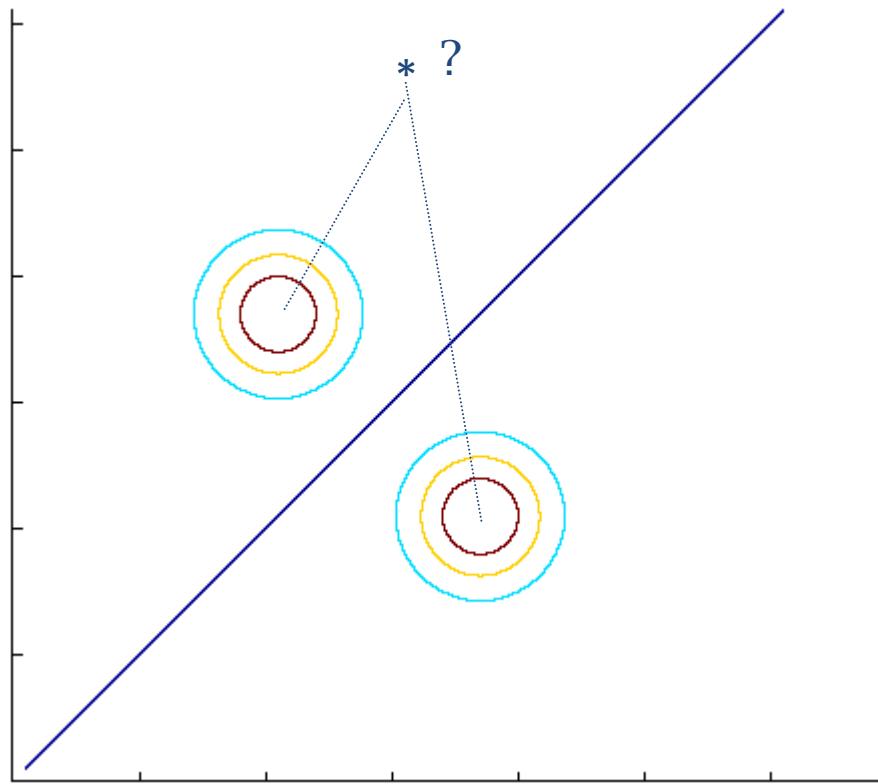
- Nearest mean classifier: Classify based on Euclidean distance to the nearest mean

$$\begin{aligned}g_i(\mathbf{x}) &= -\frac{\|\mathbf{x} - \mathbf{m}_i\|^2}{2s^2} + \log \hat{P}(C_i) \\&= -\frac{1}{2s^2} \sum_{j=1}^d (x_j - m_{ij})^2 + \log \hat{P}(C_i)\end{aligned}$$

- Each mean can be considered a prototype or template and this is template matching



Diagonal S , equal variances





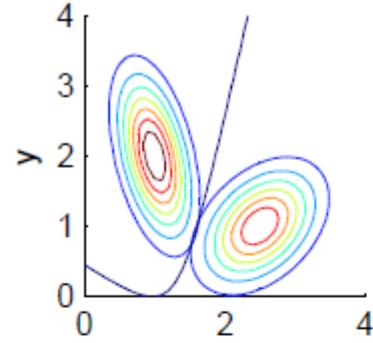
Model Selection

Assumption	Covariance matrix	No of parameters
Shared, Hyperspheric	$\mathbf{S}_i = \mathbf{S} = s^2 \mathbf{I}$	1
Shared, Axis-aligned	$\mathbf{S}_i = \mathbf{S}$, with $s_{ij} = 0$	d
Shared, Hyperellipsoidal	$\mathbf{S}_i = \mathbf{S}$	$d(d+1)/2$
Different, Hyperellipsoidal	\mathbf{S}_i	$K d(d+1)/2$

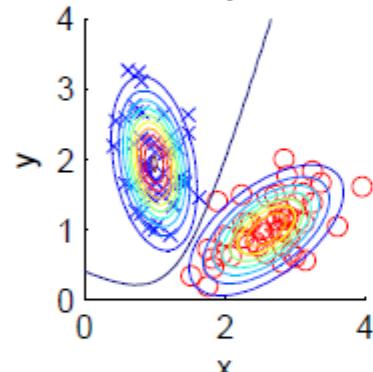
- As we increase complexity (less restricted \mathbf{S}), bias decreases and variance increases
- Assume simple models (allow some bias) to control variance (regularization)



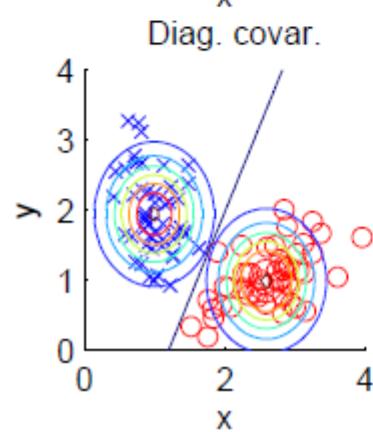
Population likelihoods and posteriors



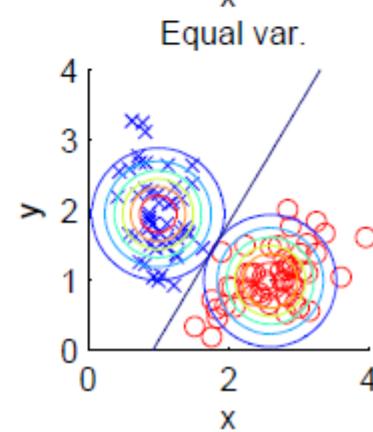
Arbitrary covar.



Shared covar.



Diag. covar.



Equal var.



Discrete Features

- Binary features: $p_{ij} \equiv p(x_j=1|C_i)$
if x_j are independent (Naive Bayes')

$$p(x|C_i) = \prod_{j=1}^d p_{ij}^{x_j} (1 - p_{ij})^{(1-x_j)}$$

the discriminant is linear

$$\begin{aligned}g_i(\mathbf{x}) &= \log p(\mathbf{x}|C_i) + \log P(C_i) \\&= \sum_j [x_j \log p_{ij} + (1 - x_j) \log (1 - p_{ij})] + \log P(C_i)\end{aligned}$$

Estimated parameters

$$\hat{p}_{ij} = \frac{\sum_t x_j^t r_i^t}{\sum_t r_i^t}$$



Discrete Features

- Multinomial (1-of- n_j) features: $x_j \in \{v_1, v_2, \dots, v_{n_j}\}$

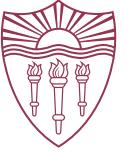
$$p_{ijk} \equiv p(z_{jk}=1 | C_i) = p(x_j=v_k | C_i)$$

if x_j are independent

$$p(\mathbf{x} | C_i) = \prod_{j=1}^d \prod_{k=1}^{n_j} p_{ijk}^{z_{jk}}$$

$$g_i(\mathbf{x}) = \sum_j \sum_k z_{jk} \log p_{ijk} + \log P(C_i)$$

$$\hat{p}_{ijk} = \frac{\sum_t z_{jk}^t r_i^t}{\sum_t r_i^t}$$



PARAMETRIC CLASSIFIERS



Likelihood- vs. Discriminant-based Classification

- **Likelihood-based**: Assume a model for $p(\mathbf{x}|C_i)$, use Bayes' rule to calculate $P(C_i|\mathbf{x})$
$$g_i(\mathbf{x}) = \log P(C_i|\mathbf{x})p(\mathbf{x})$$
- **Discriminant-based**: Assume a model for $g_i(\mathbf{x}|\Phi_i)$; no density estimation
- Estimating the boundaries is enough; no need to accurately estimate the densities inside the boundaries
- Again, this is because it is an unnormalized version of the posterior. If the goal is simply to predict the most likely class, the probabilities are not as high a concern
- **BUT**: posterior will give you confidence in each class, which is needed in many situations



Linear Discriminant

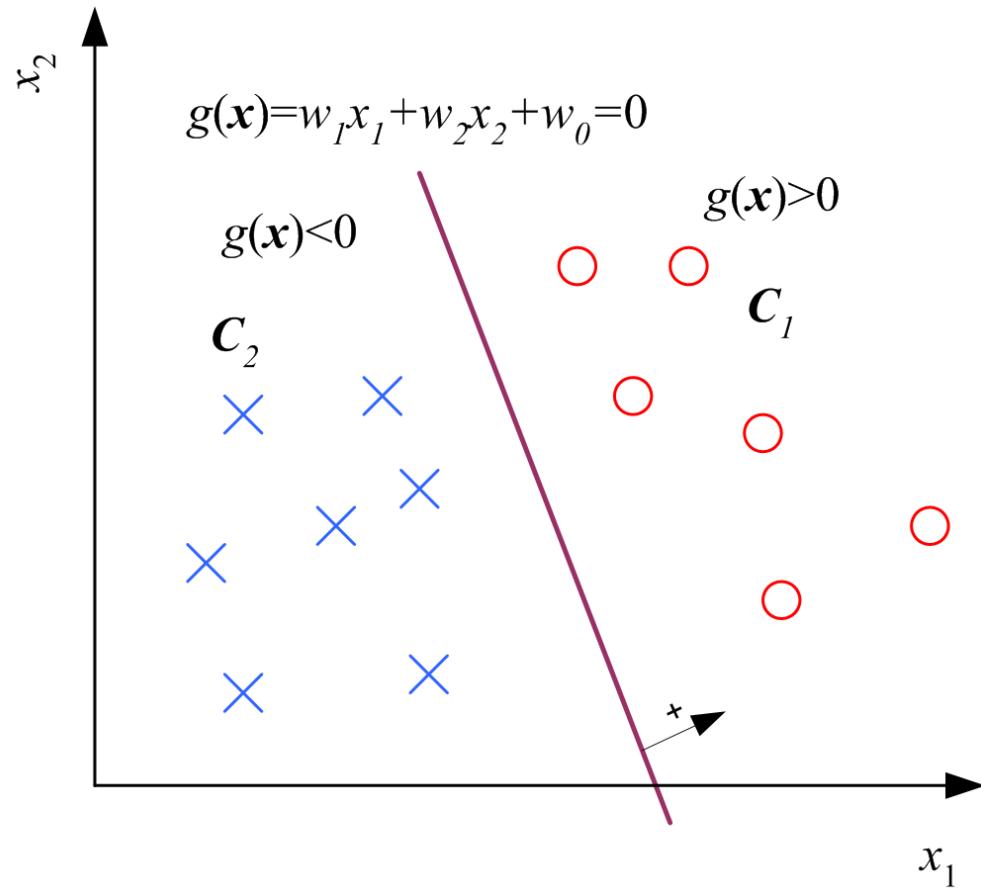
- Linear discriminant:

$$g_i(\mathbf{x} | \mathbf{w}_i, w_{i0}) = \mathbf{w}_i^T \mathbf{x} + w_{i0} = \sum_{j=1}^d w_{ij} x_j + w_{i0}$$

- Advantages:
 - Simple: $O(d)$ space/computation
 - Knowledge extraction: Weighted sum of attributes; positive/negative weights, magnitudes (credit scoring)
 - Optimal when $p(\mathbf{x} | C_i)$ are Gaussian with shared cov matrix; useful when classes are (almost) linearly separable



Two Classes

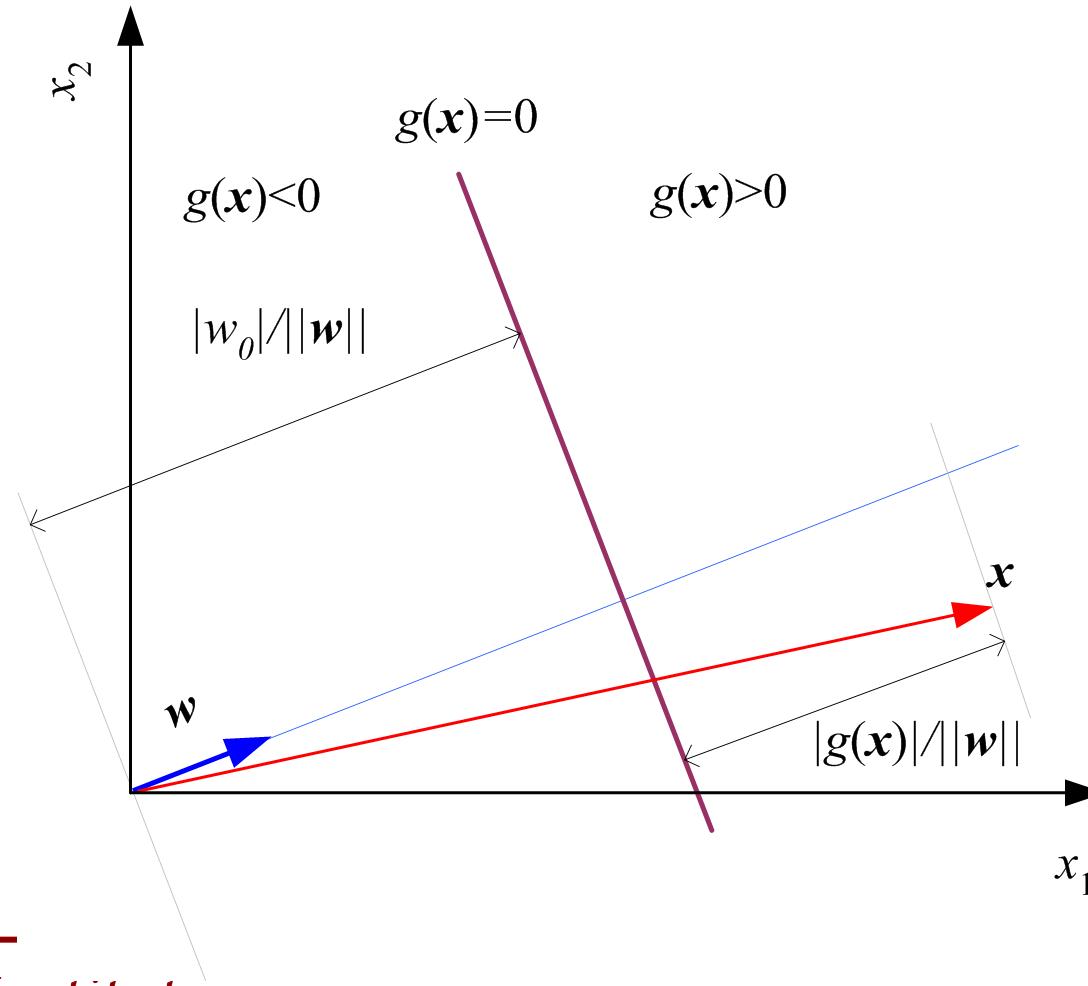


$$\begin{aligned}g(\mathbf{x}) &= g_1(\mathbf{x}) - g_2(\mathbf{x}) \\&= (\mathbf{w}_1^T \mathbf{x} + w_{10}) - (\mathbf{w}_2^T \mathbf{x} + w_{20}) \\&= (\mathbf{w}_1 - \mathbf{w}_2)^T \mathbf{x} + (w_{10} - w_{20}) \\&= \mathbf{w}^T \mathbf{x} + w_0\end{aligned}$$

choose $\begin{cases} C_1 & \text{if } g(\mathbf{x}) > 0 \\ C_2 & \text{otherwise} \end{cases}$



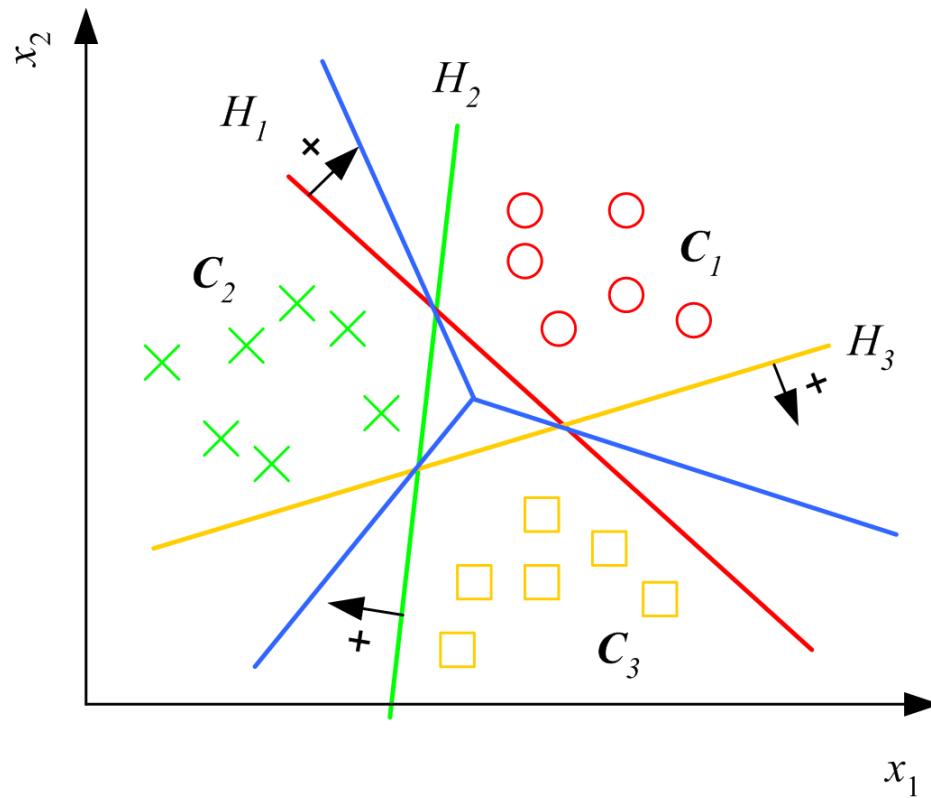
Geometry





Multiple Classes

$$g_i(\mathbf{x} | \mathbf{w}_i, w_{i0}) = \mathbf{w}_i^T \mathbf{x} + w_{i0}$$

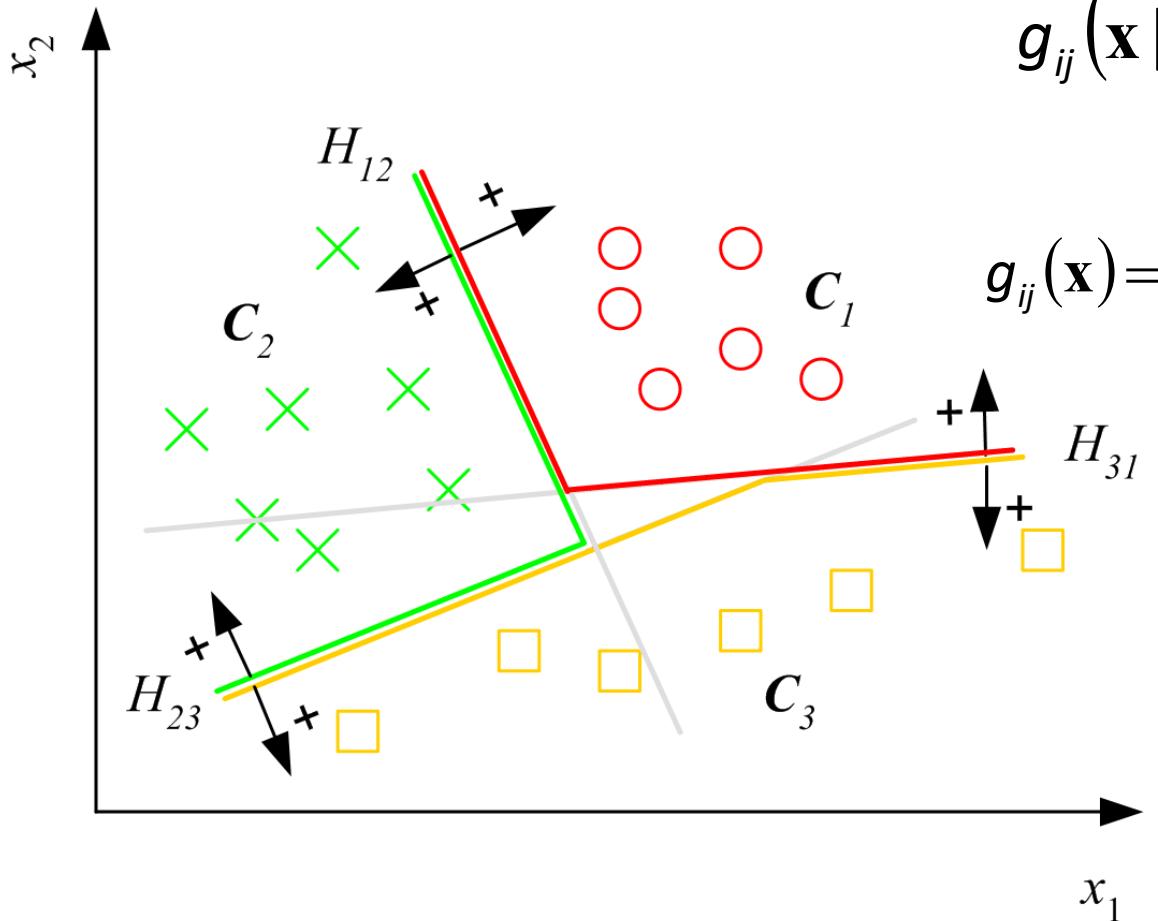


Choose C_i if

$$g_i(\mathbf{x}) = \max_{j=1}^K g_j(\mathbf{x})$$

Classes are
linearly separable

Pairwise Separation



$$g_{ij}(\mathbf{x} | \mathbf{w}_{ij}, w_{ij0}) = \mathbf{w}_{ij}^T \mathbf{x} + w_{ij0}$$

$$g_{ij}(\mathbf{x}) = \begin{cases} > 0 & \text{if } \mathbf{x} \in C_i \\ \leq 0 & \text{if } \mathbf{x} \in C_j \\ \text{don't care} & \text{otherwise} \end{cases}$$

choose C_i if
 $\forall j \neq i, g_{ij}(\mathbf{x}) > 0$



From Discriminants to Posteriors

When $p(\mathbf{x} | C_i) \sim N(\boldsymbol{\mu}_i, \Sigma)$

$$g_i(\mathbf{x} | \mathbf{w}_i, w_{i0}) = \mathbf{w}_i^T \mathbf{x} + w_{i0}$$

$$\mathbf{w}_i = \Sigma^{-1} \boldsymbol{\mu}_i \quad w_{i0} = -\frac{1}{2} \boldsymbol{\mu}_i^T \Sigma^{-1} \boldsymbol{\mu}_i + \log P(C_i)$$

$$y \equiv P(C_1 | \mathbf{x}) \text{ and } P(C_2 | \mathbf{x}) = 1 - y$$

choose C_1 if $\begin{cases} y > 0.5 \\ y/(1-y) > 1 \quad \text{and } C_2 \text{ otherwise} \\ \log [y/(1-y)] > 0 \end{cases}$



$$\begin{aligned}\text{logit}(P(C_1 | \mathbf{x})) &= \log \frac{P(C_1 | \mathbf{x})}{1 - P(C_1 | \mathbf{x})} = \log \frac{P(C_1 | \mathbf{x})}{P(C_2 | \mathbf{x})} \\ &= \log \frac{p(\mathbf{x} | C_1)}{p(\mathbf{x} | C_2)} + \log \frac{P(C_1)}{P(C_2)} \\ &= \log \frac{(2\pi)^{-d/2} |\Sigma|^{-1/2} \exp[-(1/2)(\mathbf{x} - \boldsymbol{\mu}_1)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_1)]}{(2\pi)^{-d/2} |\Sigma|^{-1/2} \exp[-(1/2)(\mathbf{x} - \boldsymbol{\mu}_2)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_2)]} + \log \frac{P(C_1)}{P(C_2)} \\ &= \mathbf{w}^T \mathbf{x} + w_0\end{aligned}$$

$$\text{where } \mathbf{w} = \Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \quad w_0 = -\frac{1}{2}(\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2)^T \Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$$

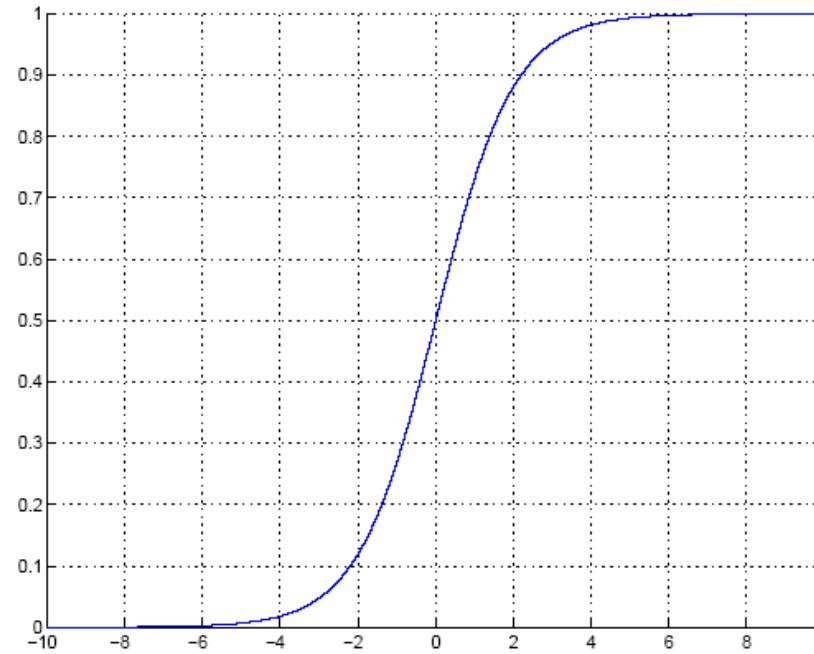
The inverse of logit

$$\log \frac{P(C_1 | \mathbf{x})}{1 - P(C_1 | \mathbf{x})} = \mathbf{w}^T \mathbf{x} + w_0$$

$$P(C_1 | \mathbf{x}) = \text{sigmoid}(\mathbf{w}^T \mathbf{x} + w_0) = \frac{1}{1 + \exp[-(\mathbf{w}^T \mathbf{x} + w_0)]}$$



Sigmoid (Logistic) Function



Calculate $g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$ and choose C_1 if $g(\mathbf{x}) > 0$, or

Calculate $y = \text{sigmoid}(\mathbf{w}^T \mathbf{x} + w_0)$ and choose C_1 if $y > 0.5$

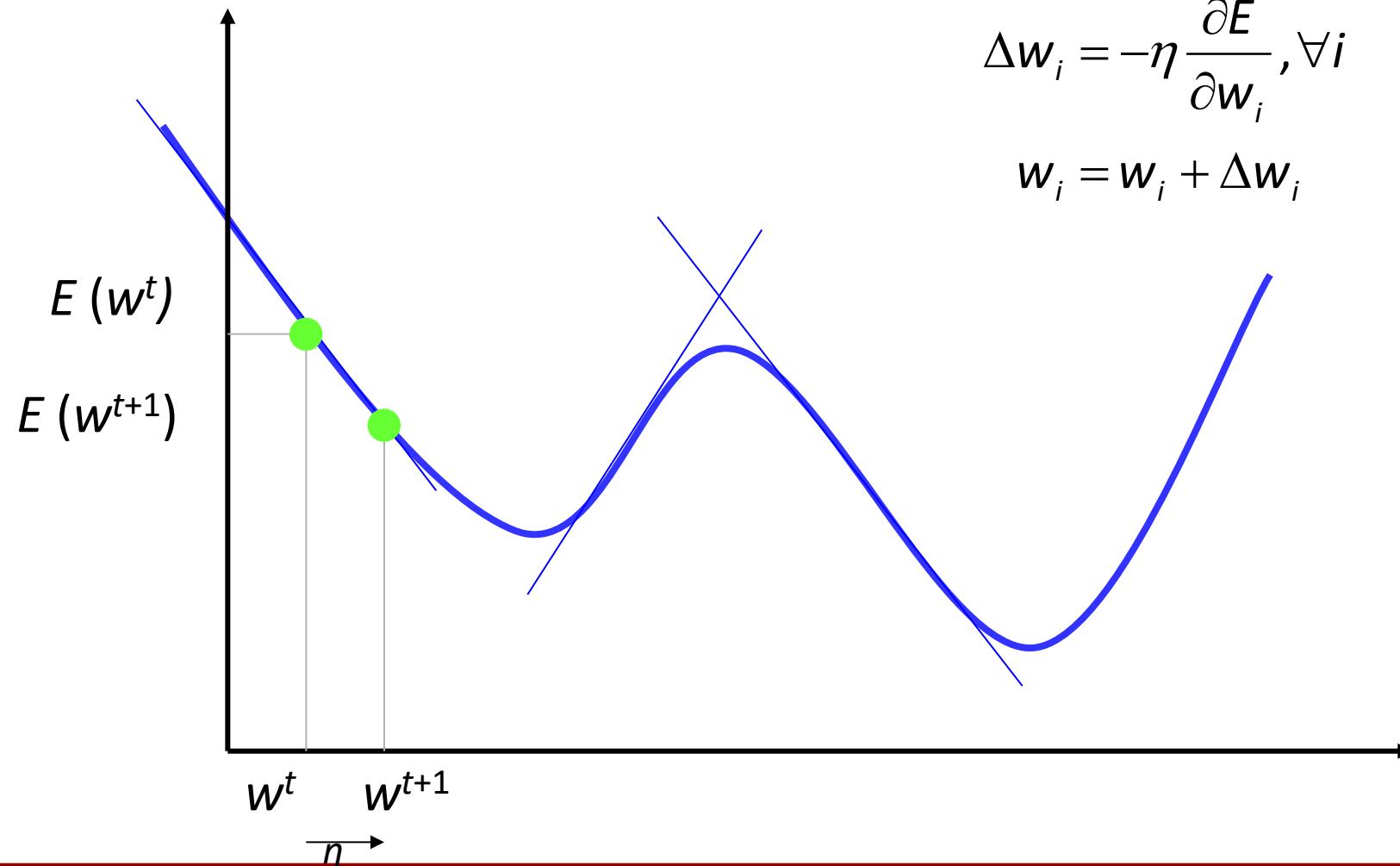


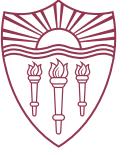
Gradient-Descent

- We have not yet discussed how to estimate the parameters...
- $E(\mathbf{w} | \mathcal{X})$ is error with parameters \mathbf{w} on sample \mathcal{X}
$$\mathbf{w}^* = \arg \min_{\mathbf{w}} E(\mathbf{w} | \mathcal{X})$$
- Gradient
$$\nabla_{\mathbf{w}} E = \left[\frac{\partial E}{\partial \mathbf{w}_1}, \frac{\partial E}{\partial \mathbf{w}_2}, \dots, \frac{\partial E}{\partial \mathbf{w}_d} \right]^T$$
- Gradient-descent:
Starts from random \mathbf{w} and updates \mathbf{w} iteratively in the negative direction of gradient



Gradient-Descent





Logistic Discrimination

Two classes: Assume log likelihood ratio is linear

$$\log \frac{p(\mathbf{x} | C_1)}{p(\mathbf{x} | C_2)} = \mathbf{w}^T \mathbf{x} + w_0^o$$

$$\begin{aligned} \text{logit}(P(C_1 | \mathbf{x})) &= \log \frac{P(C_1 | \mathbf{x})}{1 - P(C_1 | \mathbf{x})} = \log \frac{p(\mathbf{x} | C_1)}{p(\mathbf{x} | C_2)} + \log \frac{P(C_1)}{P(C_2)} \\ &= \mathbf{w}^T \mathbf{x} + w_0 \end{aligned}$$

$$\text{where } w_0 = w_0^o + \log \frac{P(C_1)}{P(C_2)}$$

$$y = \hat{P}(C_1 | \mathbf{x}) = \frac{1}{1 + \exp[-(\mathbf{w}^T \mathbf{x} + w_0)]}$$



Training: Two Classes

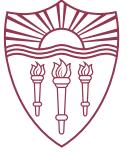
$$\mathcal{X} = \{\mathbf{x}^t, r^t\}_t \quad r^t \mid \mathbf{x}^t \sim \text{Bernoulli}(y^t)$$

$$y = P(C_1 \mid \mathbf{x}) = \frac{1}{1 + \exp[-(\mathbf{w}^\top \mathbf{x} + w_0)]}$$

$$l(\mathbf{w}, w_0 \mid \mathcal{X}) = \prod_t (y^t)^{(r^t)} (1 - y^t)^{(1 - r^t)}$$

$$E = -\log l$$

$$E(\mathbf{w}, w_0 \mid \mathcal{X}) = -\sum_t r^t \log y^t + (1 - r^t) \log (1 - y^t)$$



Training: Gradient-Descent

$$E(\mathbf{w}, w_0 | \mathcal{X}) = -\sum_t r^t \log y^t + (1 - r^t) \log (1 - y^t)$$

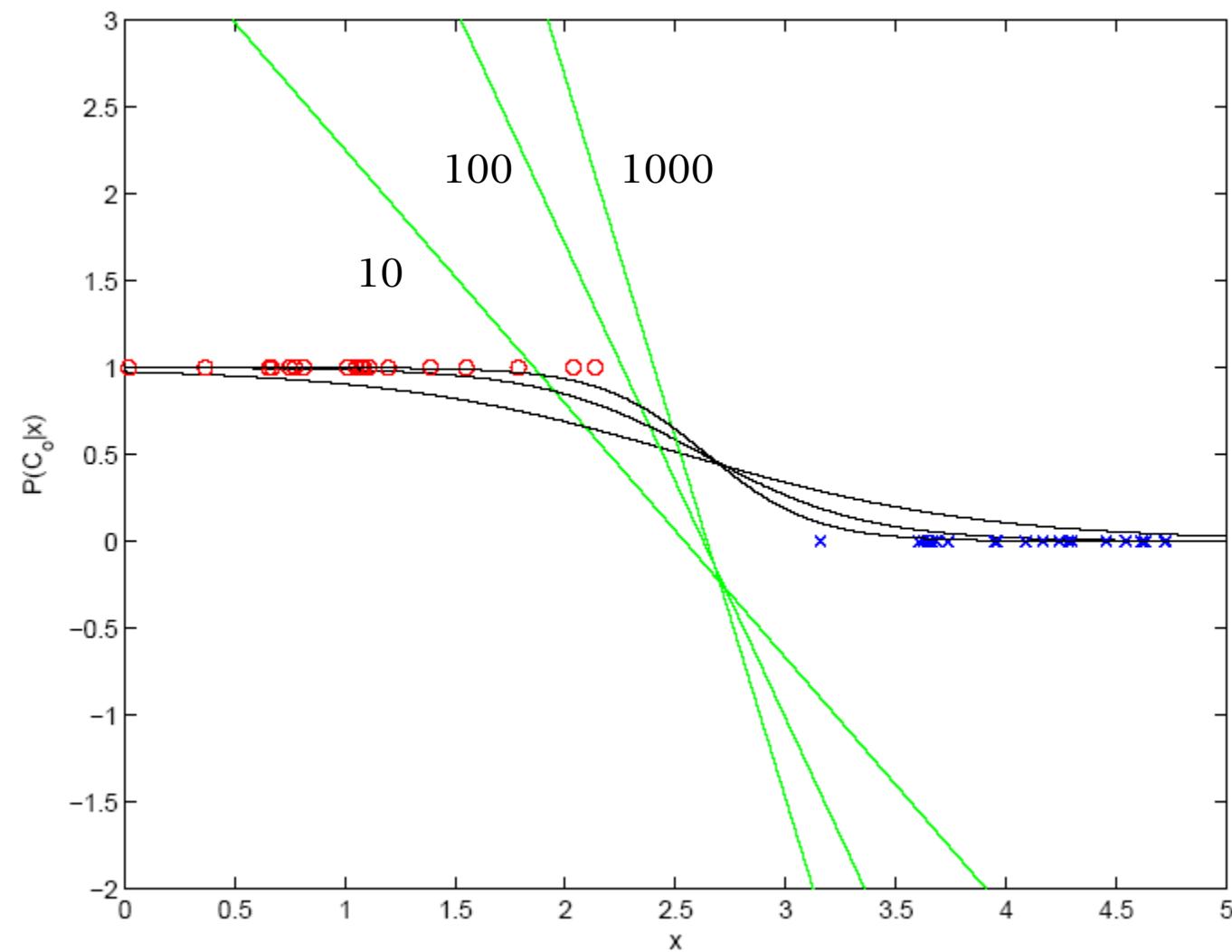
If $y = \text{sigmoid}(a)$ $\frac{dy}{da} = y(1 - y)$

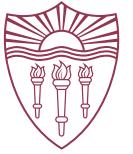
$$\begin{aligned}\Delta w_j &= -\eta \frac{\partial E}{\partial w_j} = \eta \sum_t \left(\frac{r^t}{y^t} - \frac{1 - r^t}{1 - y^t} \right) y^t (1 - y^t) x_j^t \\ &= \eta \sum_t (r^t - y^t) x_j^t, j = 1, \dots, d\end{aligned}$$

$$\Delta w_0 = -\eta \frac{\partial E}{\partial w_0} = \eta \sum_t (r^t - y^t)$$



```
For  $j = 0, \dots, d$ 
     $w_j \leftarrow \text{rand}(-0.01, 0.01)$ 
Repeat
    For  $j = 0, \dots, d$ 
         $\Delta w_j \leftarrow 0$ 
    For  $t = 1, \dots, N$ 
         $o \leftarrow 0$ 
        For  $j = 0, \dots, d$ 
             $o \leftarrow o + w_j x_j^t$ 
         $y \leftarrow \text{sigmoid}(o)$ 
         $\Delta w_j \leftarrow \Delta w_j + (r^t - y)x_j^t$ 
    For  $j = 0, \dots, d$ 
         $w_j \leftarrow w_j + \eta \Delta w_j$ 
Until convergence
```





K>2 Classes

$$\mathcal{X} = \{\mathbf{x}^t, \mathbf{r}^t\}_t \quad r^t \mid \mathbf{x}^t \sim \text{Mult}_K(1, \mathbf{y}^t)$$

$$\log \frac{p(\mathbf{x} | C_i)}{p(\mathbf{x} | C_K)} = \mathbf{w}_i^T \mathbf{x} + w_{i0}^o$$

$$y_i = \hat{P}(C_i | \mathbf{x}) = \frac{\exp[\mathbf{w}_i^T \mathbf{x} + w_{i0}^o]}{\sum_{j=1}^K \exp[\mathbf{w}_j^T \mathbf{x} + w_{j0}^o]}, i = 1, \dots, K \quad \text{softmax}$$

$$I(\{\mathbf{w}_i, w_{i0}\}_i | \mathcal{X}) = \prod_t \prod_i (y_i^t)^{(r_i^t)}$$

$$E(\{\mathbf{w}_i, w_{i0}\}_i | \mathcal{X}) = - \sum_t r_i^t \log y_i^t$$

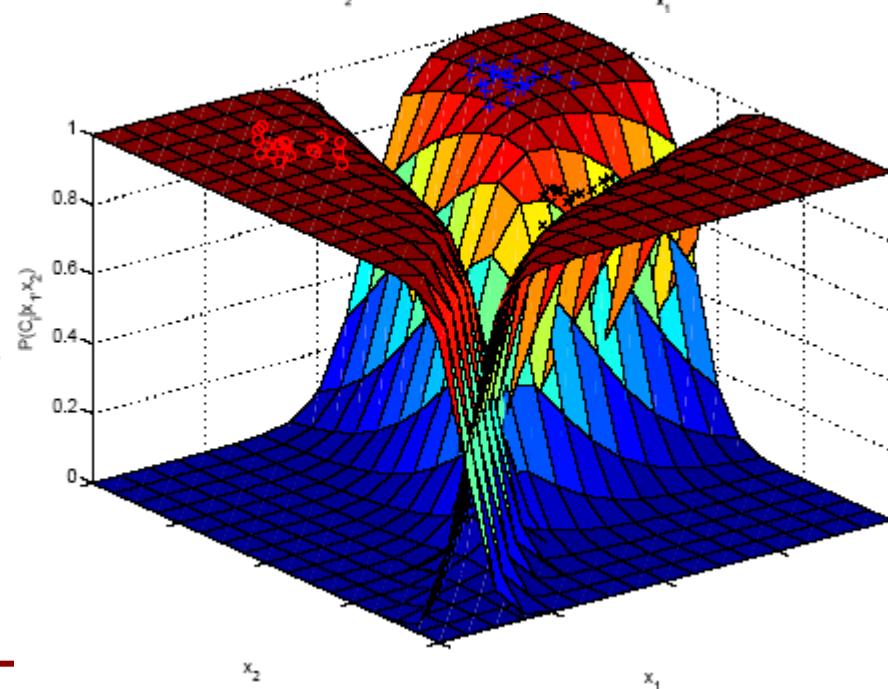
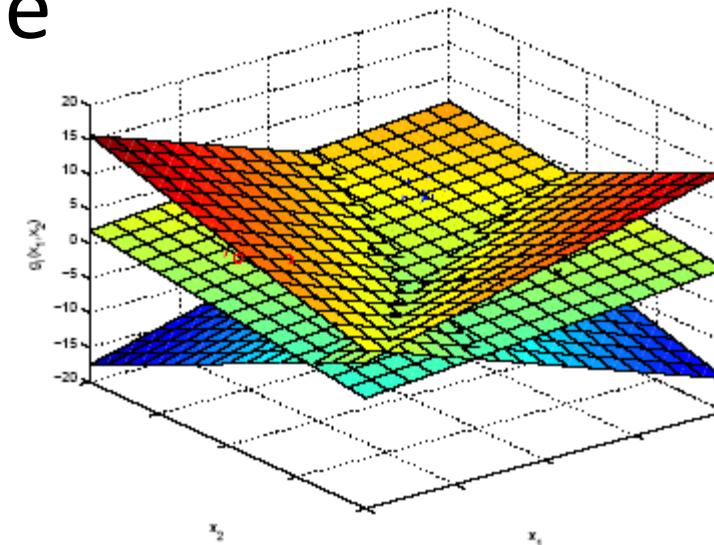
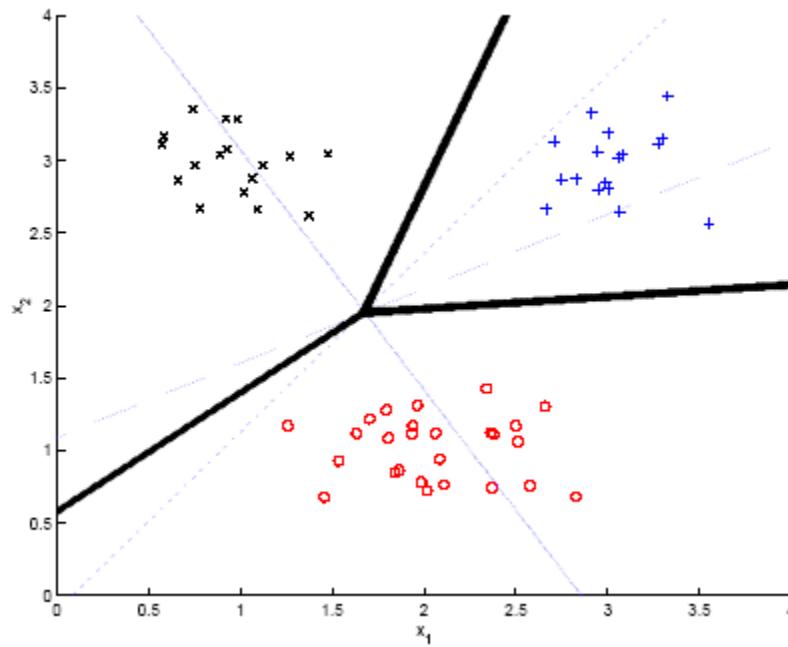
$$\Delta \mathbf{w}_j = \eta \sum_t (r_j^t - y_j^t) \mathbf{x}^t \quad \Delta w_{j0} = \eta \sum_t (r_j^t - y_j^t)$$



```
For  $i = 1, \dots, K$ , For  $j = 0, \dots, d$ ,  $w_{ij} \leftarrow \text{rand}(-0.01, 0.01)$ 
Repeat
    For  $i = 1, \dots, K$ , For  $j = 0, \dots, d$ ,  $\Delta w_{ij} \leftarrow 0$ 
    For  $t = 1, \dots, N$ 
        For  $i = 1, \dots, K$ 
             $o_i \leftarrow 0$ 
            For  $j = 0, \dots, d$ 
                 $o_i \leftarrow o_i + w_{ij}x_j^t$ 
            For  $i = 1, \dots, K$ 
                 $y_i \leftarrow \exp(o_i) / \sum_k \exp(o_k)$ 
        For  $i = 1, \dots, K$ 
            For  $j = 0, \dots, d$ 
                 $\Delta w_{ij} \leftarrow \Delta w_{ij} + (r_i^t - y_i)x_j^t$ 
    For  $i = 1, \dots, K$ 
        For  $j = 0, \dots, d$ 
             $w_{ij} \leftarrow w_{ij} + \eta \Delta w_{ij}$ 
Until convergence
```



Example





Generalizing the Linear Model

- Quadratic:

$$\log \frac{p(\mathbf{x} | C_i)}{p(\mathbf{x} | C_K)} = \mathbf{x}^T \mathbf{W}_i \mathbf{x} + \mathbf{w}_i^T \mathbf{x} + w_{i0}$$

- Sum of basis functions:

$$\log \frac{p(\mathbf{x} | C_i)}{p(\mathbf{x} | C_K)} = \mathbf{w}_i^T \phi(\mathbf{x}) + w_{i0}$$

where $\phi(\mathbf{x})$ are basis functions. Examples:

- Hidden units in neural networks (Chapters 11 and 12)
- Kernels in SVM (Chapter 13)



Discrimination by Regression

- Classes are NOT mutually exclusive and exhaustive

$$r^t = y^t + \varepsilon \text{ where } \varepsilon \sim \mathcal{N}(0, \sigma^2)$$

$$y^t = \text{sigmoid}(\mathbf{w}^\top \mathbf{x}^t + w_0) = \frac{1}{1 + \exp[-(\mathbf{w}^\top \mathbf{x}^t + w_0)]}$$

$$l(\mathbf{w}, w_0 | \mathcal{X}) = \prod_t \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(r^t - y^t)^2}{2\sigma^2}\right]$$

$$E(\mathbf{w}, w_0 | \mathcal{X}) = \frac{1}{2} \sum_t (r^t - y^t)^2$$

$$\Delta \mathbf{w} = \eta \sum_t (r^t - y^t) y^t (1 - y^t) \mathbf{x}^t$$