

Learning classifiers using generative models

1 / 29

Review: conditional probability

Let (Ω, \mathbb{P}) be a probability space.
(Ω is the sample space; \mathbb{P} is the probability distribution.)

For any events $A, B \subseteq \Omega$,

$$\mathbb{P}(A | B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)}.$$

Bayes' rule:

$$\mathbb{P}(A | B) = \frac{\mathbb{P}(A) \cdot \mathbb{P}(B | A)}{\mathbb{P}(B)}.$$

Let $E, H_0, H_1 \subseteq \Omega$. Conditioned on E , which of H_0 and H_1 is more probable?

Compare $\mathbb{P}(H_0) \cdot \mathbb{P}(E | H_0)$ to $\mathbb{P}(H_1) \cdot \mathbb{P}(E | H_1)$.

2 / 29

Conditional probability example

Suppose result of test for genetic disease is correct with probability 95%, and suppose the disease is rare: any given person has disease with probability 1%.

Question: If test comes back positive for disease, is it more likely that you have disease or do not?

E	=	test comes back positive for disease
H_0	=	do not have disease
H_1	=	have disease
$\mathbb{P}(E H_0)$	=	0.05
$\mathbb{P}(E H_1)$	=	0.95
$\mathbb{P}(H_0)$	=	0.99
$\mathbb{P}(H_1)$	=	0.01

Want to compare $\mathbb{P}(H_0 | E)$ to $\mathbb{P}(H_1 | E)$, so compare

$$\mathbb{P}(H_0) \cdot \mathbb{P}(E | H_0) = 0.99 \cdot 0.05 \quad \text{and} \quad \mathbb{P}(H_1) \cdot \mathbb{P}(E | H_1) = 0.01 \cdot 0.95.$$

3 / 29

Review: functions of random variables

Let $X: \Omega \rightarrow \mathcal{X}$ be a random variable.

"Function of a random variable":

For any function $g: \mathcal{X} \rightarrow \mathbb{R}$,

$$g(X) := g \circ X$$

is also a random variable:

$$g(X)(\omega) = g(X(\omega)).$$

Expected value:

$$\begin{aligned} \mathbb{E}(g(X)) &= \sum_{\omega \in \Omega} g(X(\omega)) \cdot \mathbb{P}(\omega) \\ &= \sum_{\gamma} \gamma \cdot \mathbb{P}(g(X) = \gamma) \\ &= \sum_x g(x) \cdot \mathbb{P}(X = x). \end{aligned}$$

4 / 29

Review: conditional expectation

Let $X: \Omega \rightarrow \mathcal{X}$ and $Y: \Omega \rightarrow \mathbb{R}$ be random variables.

Conditional expectation:
For any $x \in \mathcal{X}$ such that $\mathbb{P}(X = x) > 0$:

$$\mathbb{E}[Y \mid X = x] := \sum_y y \cdot \mathbb{P}(Y = y \mid X = x).$$

What is $\mathbb{E}[Y \mid X]$? A random variable!

$$\mathbb{E}[Y \mid X](\omega) = \mathbb{E}[Y \mid X = X(\omega)].$$

Law of total expectation:

$$\begin{aligned} \mathbb{E}[\mathbb{E}[Y \mid X]] &= \sum_{\omega \in \Omega} \mathbb{E}[Y \mid X = X(\omega)] \cdot \mathbb{P}(\omega) \\ &= \sum_x \mathbb{E}[Y \mid X = x] \cdot \mathbb{P}(X = x) = \mathbb{E}(Y). \end{aligned}$$

Bayes classifier (for binary classification)

- Probability distribution P over $\mathcal{X} \times \{0, 1\}$; let $(X, Y) \sim P$.
- Think of P as being comprised of two parts.
 1. Marginal distribution of X (a distribution over \mathcal{X}).
 2. Conditional distribution of Y given $X = x$, for each $x \in \mathcal{X}$:

$$\eta(x) := P(Y = 1 \mid X = x).$$

- The optimal classifier with smallest error rate (i.e., *Bayes classifier*) is

$$f^*(x) = \begin{cases} 0 & \text{if } \eta(x) \leq 1/2 \\ 1 & \text{if } \eta(x) > 1/2. \end{cases}$$

(Formal derivation on next slide.)

Formal derivation

Error rate of $f: \mathcal{X} \rightarrow \{0, 1\}$ can be written as

$$\text{err}_P(f) = P(f(X) \neq Y) = \mathbb{E}[\mathbb{1}\{f(X) \neq Y\}].$$

Define $g: \mathcal{X} \rightarrow \mathbb{R}$ by

$$g(x) := \mathbb{E}[\mathbb{1}\{f(X) \neq Y\} \mid X = x].$$

Then, by law of total probability,

$$\begin{aligned} g(x) &= P(Y = 0 \mid X = x) \cdot \mathbb{1}\{f(x) \neq 0\} + P(Y = 1 \mid X = x) \cdot \mathbb{1}\{f(x) \neq 1\} \\ &= (1 - \eta(x)) \cdot \mathbb{1}\{f(x) = 1\} + \eta(x) \cdot \mathbb{1}\{f(x) = 0\}. \end{aligned}$$

What should $f(x)$ be so that $g(x)$ is as small as possible?

Since $f(x) \in \{0, 1\}$, best to have

$$f(x) := \begin{cases} 0 & \text{if } \eta(x) \leq 1 - \eta(x) \\ 1 & \text{if } \eta(x) > 1 - \eta(x) \end{cases}$$

... which is the same as $f^*(x)$. □

Bayes classifier (for K -class classification)

- Probability distribution P over $\mathcal{X} \times \{1, 2, \dots, K\}$; let $(X, Y) \sim P$.
- Think of P as being comprised of two parts.
 1. Marginal distribution of X (a distribution over \mathcal{X}).
 2. Conditional distribution of Y given $X = x$, for each $x \in \mathcal{X}$.
- Bayes classifier:

$$f^*(x) = \arg \max_{y \in \{1, 2, \dots, K\}} P(Y = y \mid X = x).$$

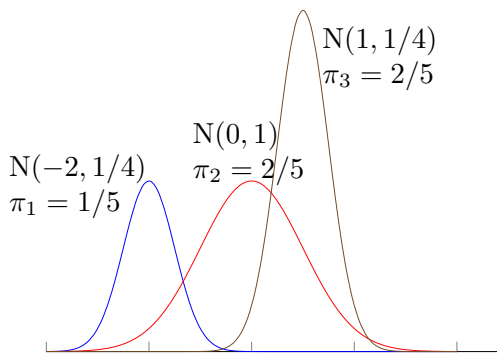
Structure of Bayes classifier

► By Bayes' rule, Bayes classifier can be written as

$$f^*(x) = \arg \max_{y \in \{1,2,\dots,K\}} P(Y = y) \cdot P(X = x \mid Y = y).$$

► Motivates thinking of P as being comprised of:

- 1. Class priors $\pi_1, \pi_2, \dots, \pi_K \in [0, 1]$, where $\pi_y = P(Y = y)$.
- 2. Class conditional distributions P_1, P_2, \dots, P_K , where P_y is conditional distribution of X given $Y = y$.



Generative models

► In context of classification problems, a *generative model* is a statistical model \mathcal{P} on $\mathcal{X} \times \{1, 2, \dots, K\}$, where each $P_\theta \in \mathcal{P}$ is

$$P_\theta(x, y) = \pi_y \cdot P_{t_y}(x),$$

where the parameters are

$$\theta = (\pi_1, \pi_2, \dots, \pi_K, t_1, t_2, \dots, t_K).$$

- Typically, all class conditional distributions $P_{t_1}, P_{t_2}, \dots, P_{t_K}$ come from same statistical model (e.g., Gaussian distribution family).
- Form of Bayes classifier corresponding to P_θ :

$$x \mapsto \arg \max_{y \in \{1,2,\dots,K\}} \pi_y \cdot P_{t_y}(x).$$

Learning a classifier

Basic approach to learning a classifier using a generative model

Suppose we observe data $D = \{(x_i, y_i)\}_{i=1}^n$, regarded as an i.i.d. sample.

0. Partition $\{x_i\}_{i=1}^n$ into D_1, D_2, \dots, D_K , where

$$D_y = \{x_i : y_i = y\}.$$

- 1. Estimate $\pi_1, \pi_2, \dots, \pi_K$ using $\{y_i\}_{i=1}^n$ (e.g., using MLE: $\hat{\pi}_y := |D_y|/n$). For each $y \in \{1, 2, \dots, K\}$: estimate t_y using D_y (e.g., using MLE).
- 2. Let classifier \hat{f} be the Bayes classifier for distribution $P_{\hat{\theta}}$ corresponding to parameter estimates

$$\hat{\theta} = (\hat{\pi}_1, \hat{\pi}_2, \dots, \hat{\pi}_K, \hat{t}_1, \hat{t}_2, \dots, \hat{t}_K),$$

i.e.,

$$\hat{f}(x) = \arg \max_{y \in \{1,2,\dots,K\}} \hat{\pi}_y \cdot P_{\hat{t}_y}(x).$$

Example: Gaussian class conditional densities

Example: $\mathcal{X} = \mathbb{R}$, $\mathcal{Y} = \{0, 1\}$, and using Gaussian class conditional densities.

Data $D = \{(x_i, y_i)\}_{i=1}^n$, regarded as an i.i.d. sample.

- 0. Split D into D_0, D_1 , where $D_0 := \{x_i : y_i = 0\}$ and $D_1 := D \setminus D_0$.
- 1. Estimate parameters:

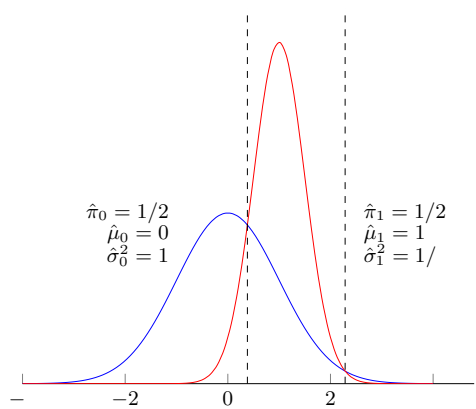
$$\begin{aligned} \hat{\pi}_0 &:= |D_0|/n, & \hat{\pi}_1 &:= |D_1|/n, \\ \hat{\mu}_0 &:= \text{sample mean}(D_0), & \hat{\mu}_1 &:= \text{sample mean}(D_1), \\ \hat{\sigma}_0^2 &:= \text{sample variance}(D_0), & \hat{\sigma}_1^2 &:= \text{sample variance}(D_1). \end{aligned}$$

2. Form classifier \hat{f} using

$$\hat{f}(x) = \arg \max_{y \in \{0,1\}} \hat{\pi}_y \cdot \varphi_{\hat{\mu}_y, \hat{\sigma}_y^2}(x)$$

where φ_{μ, σ^2} is the $N(\mu, \sigma^2)$ density.

Example: Gaussian class conditional densities



Classifier:

$$\hat{f}(x) = \begin{cases} 1 & \text{if } x \in [0.38, 2.29]; \\ 0 & \text{otherwise.} \end{cases}$$

Dotted lines = *decision boundary*.

Bayes classifiers



Suppose $\{(x_i, y_i)\}_{i=1}^n$ is an iid sample from P , and let \mathcal{P} be a generative model with parameter space Θ .

- ▶ Let $\hat{\theta} \in \Theta$ be parameter estimate obtained using data $\{(x_i, y_i)\}_{i=1}^n$.
- ▶ Let \hat{f} be the Bayes classifier for distribution $P_{\hat{\theta}} \in \mathcal{P}$.

Let f^* be the Bayes classifier for P .

- ▶ Is \hat{f} the same as f^* ?

High-dimensional feature spaces

Suppose $\mathcal{X} = \{0, 1\}^d$ or $\mathcal{X} = \mathbb{Z}_+^d$ or $\mathcal{X} = \mathbb{R}^d$ where $d > 1$.

What statistical models can we use for the **class conditional distributions**?

- ▶ *non-parametric models*: very general, but quality may be poor for large d .
- ▶ Often have prior knowledge about *statistical dependencies* between variables. Leverage this knowledge to form a *graphical model*.
- ▶ Some simple models: *multivariate Gaussians*, *product distributions*.

Product distributions on $\{0, 1\}^d$

Suppose $\mathcal{X} = \{0, 1\}^d$, and let \mathcal{P} be all *product distributions* on $\{0, 1\}^d$.

Parameters $\mu = (\mu_1, \mu_2, \dots, \mu_d) \in [0, 1]^d$:

$$P_{\mu}(x) = \prod_{j=1}^d \mu_j^{x_j} (1 - \mu_j)^{1-x_j} \quad \text{for all } x = (x_1, x_2, \dots, x_d) \in \{0, 1\}^d.$$

If **random vector** $X = (X_1, X_2, \dots, X_d) \sim P_{\mu}$, then X_1, X_2, \dots, X_d are independent random variables, and

$$P_{\mu}(X_j = 1) = \mu_j.$$

Naïve Bayes classifiers

Generative models that use **product distributions as class conditionals**
→ *Naïve Bayes classifiers*.

(Using product distributions on $\{0, 1\}^d$ as in previous slide.)
What is the form of Bayes classifier corresponding to distribution with parameters $\theta = (\pi_1, \pi_2, \dots, \pi_K, \mu_1, \mu_2, \dots, \mu_K)$?

$$\begin{aligned} \mathbf{x} &\mapsto \arg \max_{y \in \{1, 2, \dots, K\}} \pi_y \cdot P_{\mu_y}(\mathbf{x}) \\ &= \arg \max_{y \in \{1, 2, \dots, K\}} \log(\pi_y \cdot P_{\mu_y}(\mathbf{x})) \\ &= \arg \max_{y \in \{1, 2, \dots, K\}} \log \left(\pi_y \cdot \prod_{j=1}^d \mu_{y,j}^{x_j} (1 - \mu_{y,j})^{1-x_j} \right) \\ &= \arg \max_{y \in \{1, 2, \dots, K\}} \log \pi_y + \sum_{j=1}^d x_j \log \mu_{y,j} + (1 - x_j) \log(1 - \mu_{y,j}) \\ &= \arg \max_{y \in \{1, 2, \dots, K\}} b_y + \langle \mathbf{w}_y, \mathbf{x} \rangle \end{aligned}$$

for some appropriate definition of $b_y \in \mathbb{R}$ and $\mathbf{w}_y \in \mathbb{R}^d$ in terms of π_y and μ_y .

17 / 29

Standard Gaussian distributions on \mathbb{R}^d

Standard normal (Gaussian) distribution on \mathbb{R}^1

$X \sim N(0, 1)$, density

$$\varphi_{0,1}(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right) \text{ for all } x \in \mathbb{R}.$$

Standard normal (Gaussian) distribution on \mathbb{R}^d

$\mathbf{X} = (X_1, X_2, \dots, X_d) \sim N(\mathbf{0}, \mathbf{I})$, density

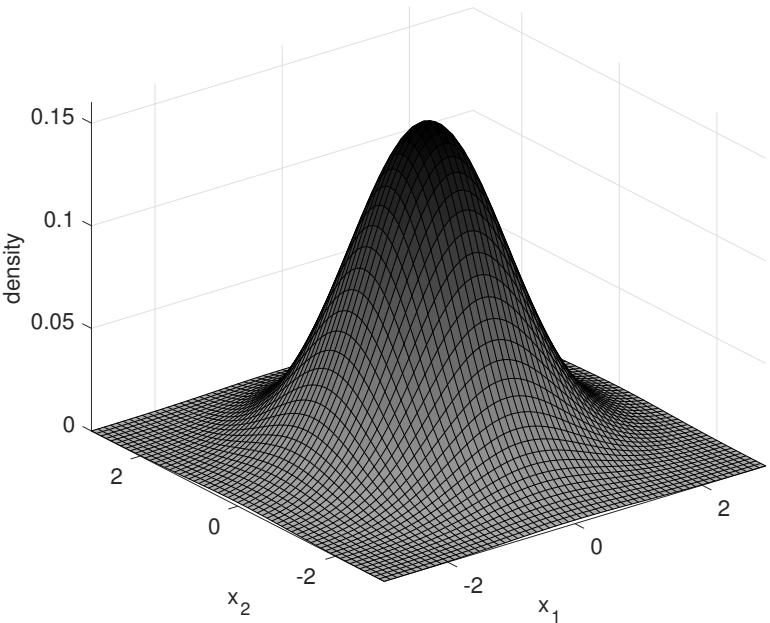
$$\varphi_{\mathbf{0},\mathbf{I}}(\mathbf{x}) = \prod_{i=1}^d \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x_i^2}{2}\right) \text{ for all } \mathbf{x} = (x_1, x_2, \dots, x_d) \in \mathbb{R}^d.$$

Usually written as

$$\varphi_{\mathbf{0},\mathbf{I}}(\mathbf{x}) = \frac{1}{(2\pi)^{d/2}} \exp\left(-\frac{\|\mathbf{x}\|_2^2}{2}\right).$$

18 / 29

Standard normal (Gaussian) distribution on \mathbb{R}^d



19 / 29

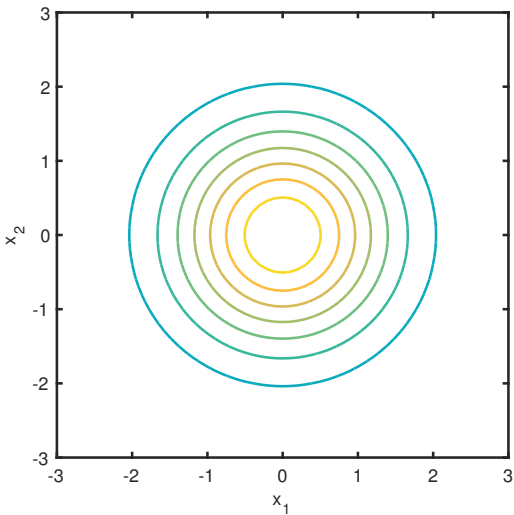
Standard normal (Gaussian) distribution on \mathbb{R}^d

- ▶ $\mathbb{E}(X_i) = 0$
- ▶ $\text{var}(X_i) = \text{cov}(X_i, X_i) = 1$
- ▶ $\text{cov}(X_i, X_j) = 0$ for $i \neq j$.

Arrange means into a vector and covariances in a $d \times d$ matrix:

$$\mathbb{E}(\mathbf{X}) = \mathbf{0}, \quad \text{cov}(\mathbf{X}) = \mathbf{I}$$

(zero vector and identity matrix).



Contours of equal standard normal density in \mathbb{R}^2

20 / 29

General Gaussian distributions on \mathbb{R}^d

(General) Gaussian distributions on \mathbb{R}^d come from applying two operations to another (e.g., the standard) Gaussian distribution:

$$\overbrace{x \mapsto Ax}^{\text{linear map}} \mapsto \underbrace{Ax + \mu}_{\text{translation}}$$

for some vector $\mu \in \mathbb{R}^d$ and invertible linear map $A \in \mathbb{R}^{d \times d}$.

Fact: Let $\mu \in \mathbb{R}^d$ be any vector, and $A \in \mathbb{R}^{d \times d}$ be any invertible matrix. For any random vector X in \mathbb{R}^d , the random vector $Y = AX + \mu$ satisfies

$$\mathbb{E}(Y) = \mu, \quad \text{cov}(Y) = AA^\top.$$

Furthermore, if $X \sim N(0, I)$, then $Y \sim N(\mu, AA^\top)$.

Examples of linear maps

Write $x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$.

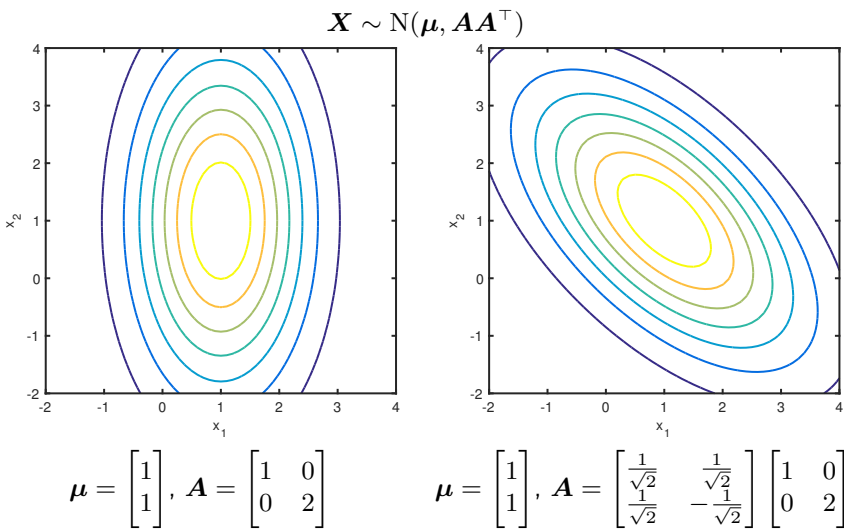
1. If $A = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}$, then $Ax = \begin{bmatrix} x_1 \\ 2x_2 \end{bmatrix}$.

(Scale coordinates x_1 and x_2 by, respectively, 1 and 2.)

2. If $A = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}$, then $Ax = x_1 \begin{bmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix} + 2x_2 \begin{bmatrix} \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} \end{bmatrix}$.

(Coordinate scaling as above, followed by rotation.)

General Gaussian distributions on \mathbb{R}^d



Gaussian parametric model

Gaussian parametric model:

$\mathcal{P} = \{P_\theta : \theta \in \Theta\}$, where $\theta = (\mu, \Sigma) \in \mathbb{R}^d \times \mathbb{S}_{++}$, with

\mathbb{S}_{++} = all symmetric positive definite (p.d.) $d \times d$ real matrices.

Recall: a matrix M is p.d. if and only if

$$v^\top M v > 0 \quad \text{for all } v \neq 0.$$

Example: Let $A \in \mathbb{R}^{d \times d}$ be invertible. For any $v \neq 0$,

$$\begin{aligned} v^\top (AA^\top) v &= (A^\top v)^\top (A^\top v) \\ &= w^\top w \quad (\text{letting } w := (A^\top)^{-1} v) \\ &= \|w\|_2^2 \geq 0. \end{aligned}$$

The only vector with zero length is 0 , so $\|w\|_2^2 = 0$ if and only if $w = 0$. Furthermore, $w = 0$ if and only if $v = A^\top w = 0$. But $v \neq 0$, so $w \neq 0$ and $\|w\|_2^2 > 0$. Hence, AA^\top is p.d.

MLE for Gaussian parameters

- Suppose we observe $\{x_i\}_{i=1}^n$ (regarded as an i.i.d. sample). What is the MLE for the Gaussian parameters (μ, Σ) ?

- Log-likelihood of $\theta = (\mu, \Sigma)$ given $\{x_i\}_{i=1}^n$:

$$\log \mathcal{L}(\theta; \{x_i\}_{i=1}^n) = \sum_{i=1}^n \left[\log \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp \left(-\frac{1}{2} (x_i - \mu)^\top \Sigma^{-1} (x_i - \mu) \right) \right].$$

- Gradient of above with respect to μ :

$$-\sum_{i=1}^n \Sigma^{-1} (x_i - \mu).$$

- Lo and behold, MLE for μ is the *sample mean*:

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i.$$

- Using *matrix derivatives* gives the MLE for Σ , the *sample covariance*:

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})(x_i - \hat{\mu})^\top$$

(where $\hat{\mu}$ is the sample mean). **What could go wrong?**

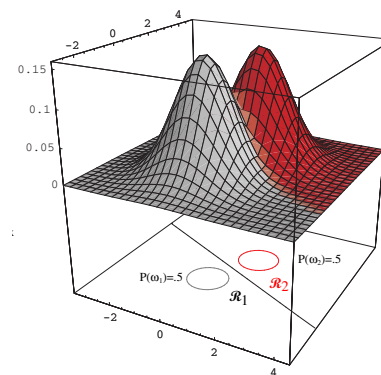
25 / 29

Multivariate Gaussian class conditionals

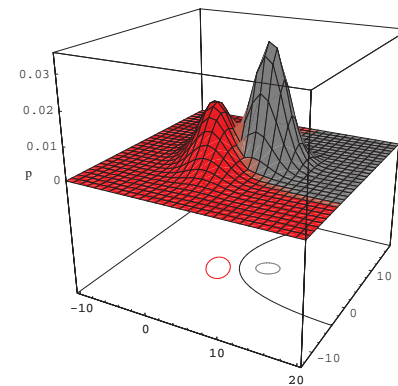
Example: $\mathcal{X} = \mathbb{R}^d$, $\mathcal{Y} = \{0, 1\}$, and using **multivariate Gaussian class conditional densities**.

Bayes classifier corresponding to distribution with parameters

$\theta = (\pi_0, \pi_1, \mu_0, \Sigma_0, \mu_1, \Sigma_1)$:



$\Sigma_0 = \Sigma_1$
Bayes classifier:
linear decision boundary



$\Sigma_0 \neq \Sigma_1$
Bayes classifier:
quadratic decision boundary

26 / 29

Generative models with parameter tying

Sometimes, we must estimate (some of) the parameters of the **class conditional distributions jointly**, rather than separately.

Example: **multivariate Gaussian class conditionals with shared covariance:**

$$t_1 = (\mu_1, \Sigma), \quad t_2 = (\mu_2, \Sigma), \quad \dots, \quad t_K = (\mu_K, \Sigma).$$

This is called *parameter tying*.

MLEs for $\pi_1, \pi_2, \dots, \pi_K, \mu_1, \mu_2, \dots, \mu_K$ given $\{(x_i, y_i)\}_{i=1}^n$:

$$\hat{\pi}_y := |D_y|/n, \quad \hat{\mu}_y := \text{sample mean}(D_y).$$

But what's the MLE for the shared class conditional covariance Σ ?

$$\hat{\Sigma} := \frac{1}{n} \sum_{y=1}^K \sum_{x_i \in D_y} (x_i - \hat{\mu}_y)(x_i - \hat{\mu}_y)^\top.$$

27 / 29

Final remarks

Some redeeming qualities of classifiers based on generative models:

- Simple recipe, many variations.
- Can leverage domain knowledge about class conditional distributions.
- Can be very efficient when K is large.

Critical drawbacks:

- Classifier relies on formula (via Bayes' rule) that assumes the estimated class priors and conditional distributions are perfect, which is not true.
- Modeling P away from decision boundary between classes is wasted effort: not necessary for good classification.

28 / 29

Key takeaways

1. Generative structure of Bayes classifier.
2. Basic properties of multivariate Gaussians.
3. Basic recipe for learning a classifier based on a generative model, and concept of parameter tying.
4. Specific generative models with product distributions and multivariate Gaussians as class conditionals.
5. High-level advantages and disadvantages of classifiers based on generative models.