

Machine Learning

4771

Instructor: Tony Jebara

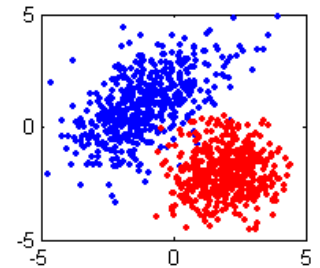
Topic 10

- Classification with Gaussians
- Regression with Gaussians
- Principal Components Analysis

Classification with Gaussians

- Have two classes, each with their own Gaussian:

$$\left\{ (x_1, y_1), \dots, (x_N, y_N) \right\} \quad x \in R^D \quad y \in \{0, 1\}$$



- Given parameters $\theta = \{\alpha, \mu_0, \Sigma_0, \mu_1, \Sigma_1\}$ we can generate iid data from $p(x, y | \theta) = p(y | \theta) p(x | y, \theta)$ by:

1) flipping a coin to get y via Bernoulli $p(y | \theta) = \alpha^y (1 - \alpha)^{1-y}$

2) sampling an x from y 'th Gaussian $p(x | y, \theta) = N(x | \mu_y, \Sigma_y)$

- Or, recover parameters from data using maximum likelihood

$$\begin{aligned} l(\theta) &= \log p(\text{data} | \theta) = \sum_{i=1}^N \log p(x_i, y_i | \theta) \\ &= \sum_{i=1}^N \log p(y_i | \theta) + \sum_{i=1}^N \log p(x_i | y_i, \theta) \\ &= \sum_{i=1}^N \log p(y_i | \alpha) + \sum_{y_i \in 0} \log p(x_i | \mu_0, \Sigma_0) + \sum_{y_i \in 1} \log p(x_i | \mu_1, \Sigma_1) \end{aligned}$$

Classification with Gaussians

- Max Likelihood can be done separately for the 3 terms

$$l = \sum_{i=1}^N \log p(y_i | \alpha) + \sum_{y_i \in 0} \log p(x_i | \mu_0, \Sigma_0) + \sum_{y_i \in 1} \log p(x_i | \mu_1, \Sigma_1)$$

- Count # of pos & neg examples (class prior): $\alpha = \frac{N_1}{N_0 + N_1}$
- Get mean & cov of negatives and mean & cov of positives:

$$\begin{aligned} \mu_0 &= \frac{1}{N_0} \sum_{y_i \in 0} x_i & \Sigma_0 &= \frac{1}{N_0} \sum_{y_i \in 0} (x_i - \mu_0)(x_i - \mu_0)^T \\ \mu_1 &= \frac{1}{N_1} \sum_{y_i \in 1} x_i & \Sigma_1 &= \frac{1}{N_1} \sum_{y_i \in 1} (x_i - \mu_1)(x_i - \mu_1)^T \end{aligned}$$

- Given (x,y) pair, can now compute likelihood $p(x, y)$
- To make classification, a bit of Decision Theory
- Without x, can compute prior guess for y $p(y)$
- Give me x, want y, I need posterior $p(y | x)$
- Bayes Optimal Decision: $\hat{y} = \arg \max_{y \in \{0,1\}} p(y | x)$
- Optimal iff we have true probability

Posterior gives Logistic

- Bayes Optimal Decision: $\hat{y} = \arg \max_{y=\{0,1\}} p(y | x)$

- To get conditional:

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

- Check which is greater: $p(y=0 | x) \geq ? \leq p(y=1 | x)$

- Or check if this is > 0.5 $p(y=1 | x) = \frac{p(x, y=1)}{p(x, y=0) + p(x, y=1)}$

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

$$= \frac{1}{\exp\left(-\log \frac{p(x, y=1)}{p(x, y=0)}\right) + 1}$$

- Get logistic squashing function of log-ratio of probability models

$$= \text{sigmoid}\left(\log \frac{p(x, y=1)}{p(x, y=0)}\right)$$

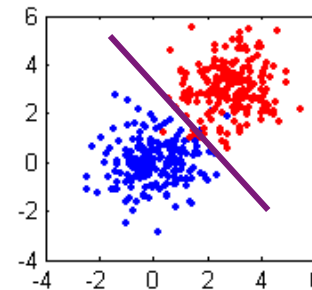
Linear or Quadratic Decisions

- Example cases, plotting decision boundary when $\alpha = 0.5$

$$\begin{aligned} p(y = 1 | x) &= \frac{p(x, y = 1)}{p(x, y = 0) + p(x, y = 1)} \\ &= \frac{\alpha N(x | \mu_1, \Sigma_1)}{(1 - \alpha) N(x | \mu_0, \Sigma_0) + \alpha N(x | \mu_1, \Sigma_1)} \end{aligned}$$

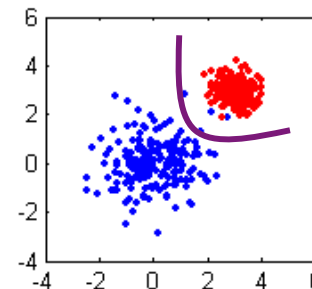
- If covariances are equal:

linear decision



- If covariances are different:

quadratic decision



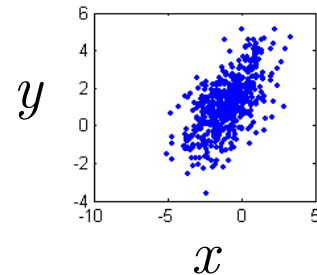
Regression with Gaussians

- Have input and output, each Gaussian:

$$\left\{ (x_1, y_1), \dots, (x_N, y_N) \right\} \quad x \in \mathbb{R}^{D_x} \quad y \in \mathbb{R}^{D_y}$$

concatenate $z_i = \begin{bmatrix} x_i \\ y_i \end{bmatrix}$

$$p(z \mid \mu, \Sigma) = \frac{1}{(2\pi)^{D/2} \sqrt{|\Sigma|}} \exp \left(-\frac{1}{2} (z - \mu)^T \Sigma^{-1} (z - \mu) \right)$$



- Maximum Likelihood is as usual for a multivariate Gaussian

$$\mu = \frac{1}{N} \sum_{i=1}^N z_i \quad \Sigma = \frac{1}{N} \sum_{i=1}^N (z_i - \mu)(z_i - \mu)^T$$

- Bayes optimal decision:

$$\hat{y} = \arg \max_{y \in \mathbb{R}} p(y \mid x)$$

- Or we can use:

$$\hat{y} = E_{p(y|x)} \{y\}$$

- Have joint, need conditional:

$$p(y \mid x) = \frac{p(x, y)}{p(x)} = \frac{p(x, y)}{\int_y p(x, y)}$$

Gaussian Marginals/Conditionals

•Conditional & marginal from joint: $p(y | x) = \frac{p(x, y)}{p(x)} = \frac{p(x, y)}{\int_y p(x, y)}$

•Conditioning the Gaussian:

$$p(z | \mu, \Sigma) = \frac{1}{(2\pi)^{D/2} \sqrt{|\Sigma|}} \exp\left(-\frac{1}{2}(z - \mu)^T \Sigma^{-1}(z - \mu)\right)$$

$$p(x, y) = \frac{1}{(2\pi)^{D/2} \sqrt{|\Sigma|}} \exp\left(-\frac{1}{2}\left(\begin{bmatrix} x \\ y \end{bmatrix} - \begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}\right)^T \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix}^{-1} \left(\begin{bmatrix} x \\ y \end{bmatrix} - \begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}\right)\right)$$

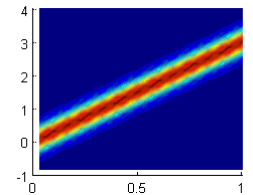
$$p(x) = \frac{1}{(2\pi)^{D_x/2} \sqrt{|\Sigma_{xx}|}} \exp\left(-\frac{1}{2}(x - \mu_x)^T \Sigma_{xx}^{-1}(x - \mu_x)\right)$$

$$= N(x | \mu_x, \Sigma_{xx})$$

$$p(y | x) = N\left(y | \mu_y + \Sigma_{yx} \Sigma_{xx}^{-1}(x - \mu_x), \Sigma_{yy} - \Sigma_{yx} \Sigma_{xx}^{-1} \Sigma_{xy}\right)$$

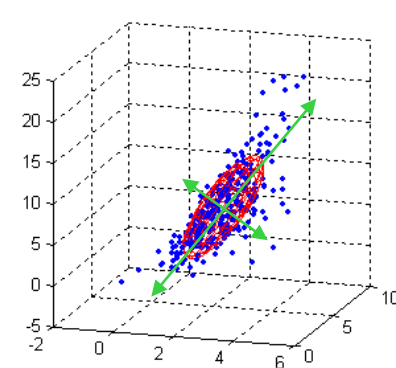
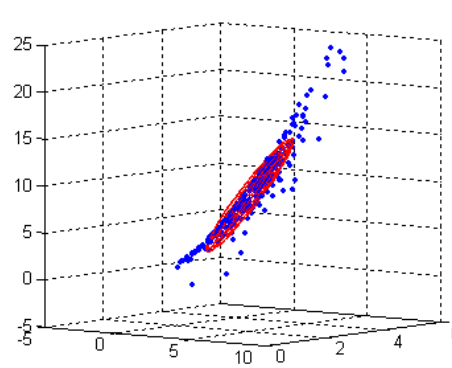
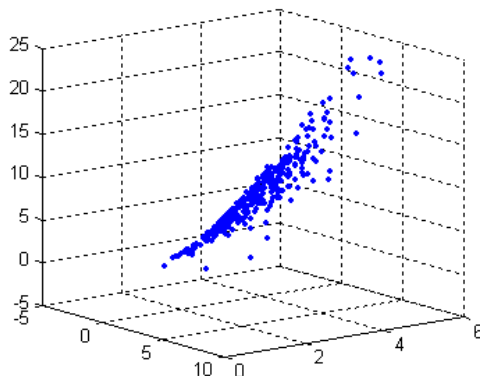
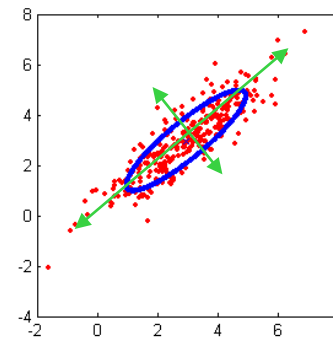
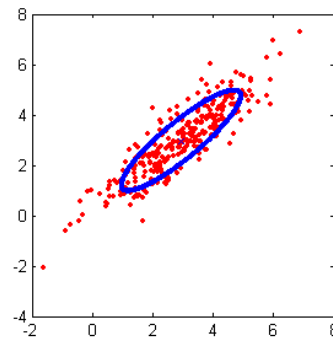
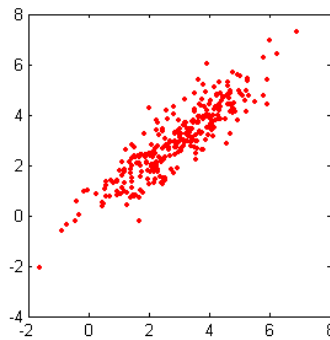
•Here argmax is expectation
which is conditional mean:

$$\hat{y} = \mu_y + \Sigma_{yx} \Sigma_{xx}^{-1}(x - \mu_x)$$



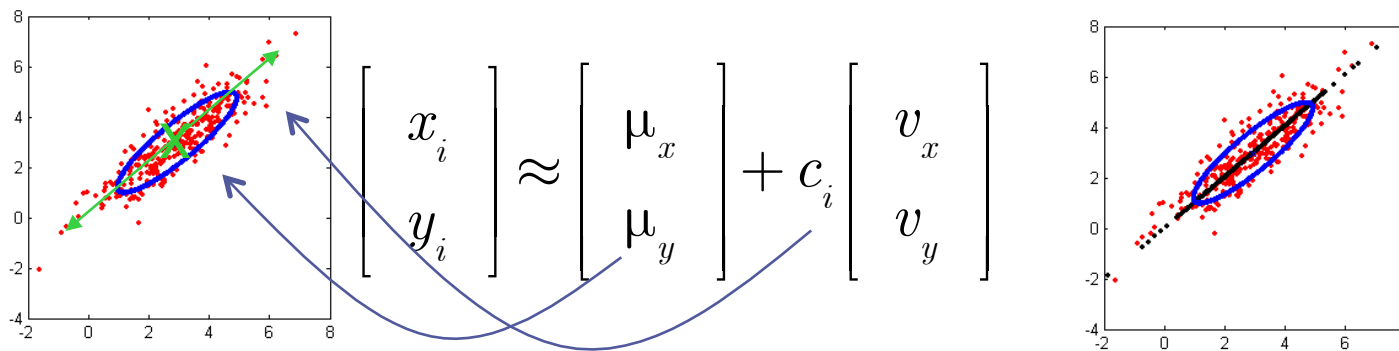
Principal Components Analysis

- Gaussians: for Classification, Regression... & Compression!
- Data can be constant in some directions, changes in others
- Use Gaussian to find directions of high/low variance



Principal Components Analysis

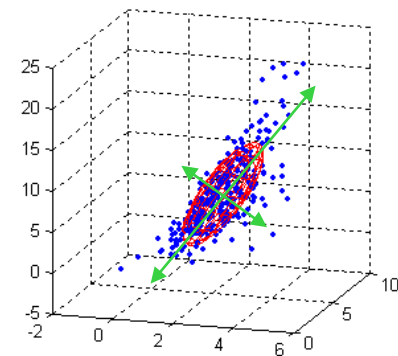
- Idea: instead of writing data in all its dimensions, only write it as mean + steps along one direction



- More generally, keep a subset of dimensions C from D (i.e. 2 of 3)

$$\vec{x}_i \approx \vec{\mu} + \sum_{j=1}^C c_{ij} \vec{v}_j$$

- Compression method: $\vec{x}_i \gg \vec{c}_i$
- Optimal directions: along eigenvectors of covariance
- Which directions to keep: highest eigenvalues (variances)



Principal Components Analysis

- If we have eigenvectors, mean and coefficients:

$$\vec{x}_i \approx \vec{\mu} + \sum_{j=1}^C c_{ij} \vec{v}_j$$

- Get eigenvectors (use eig() in Matlab): $\Sigma = V \Lambda V^T$

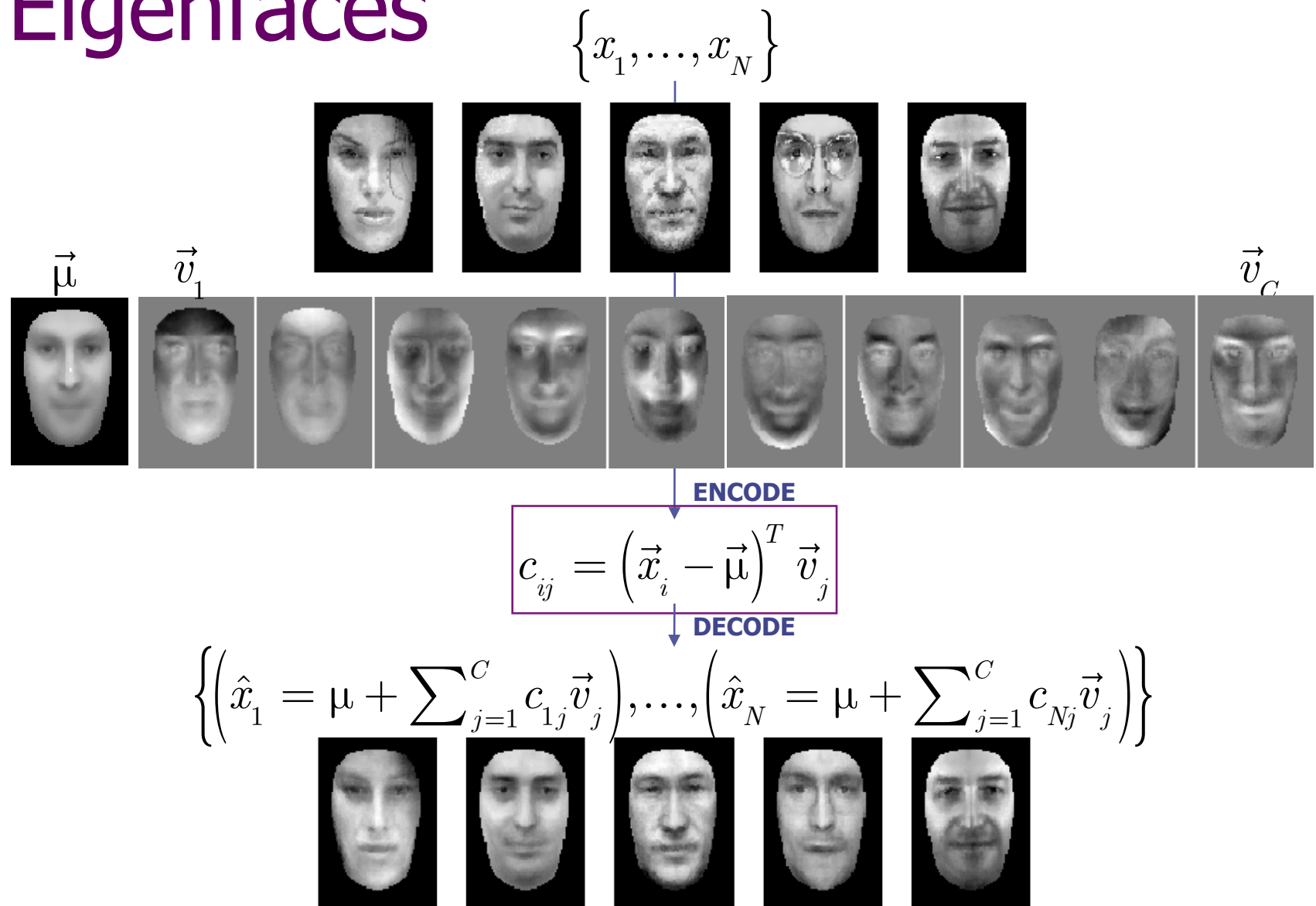
$$\begin{bmatrix} \Sigma(1,1) & \Sigma(1,2) & \Sigma(1,3) \\ \Sigma(1,2) & \Sigma(2,2) & \Sigma(2,3) \\ \Sigma(1,3) & \Sigma(2,3) & \Sigma(3,3) \end{bmatrix} = \begin{bmatrix} [\vec{v}_1] & [\vec{v}_2] & [\vec{v}_3] \end{bmatrix} \begin{bmatrix} \lambda_1 & 0 & 0 \\ 0 & \lambda_2 & 0 \\ 0 & 0 & \lambda_3 \end{bmatrix} \begin{bmatrix} [\vec{v}_1] & [\vec{v}_2] & [\vec{v}_3] \end{bmatrix}^T$$

- Eigenvectors are orthonormal: $\vec{v}_i^T \vec{v}_j = \delta_{ij}$
- In coordinates of v , Gaussian is diagonal, $\text{cov} = \Lambda$
- All eigenvalues are non-negative $\lambda_i \geq 0$
- Higher eigenvalues are higher variance, use the top C ones

$$\lambda_1 \geq \lambda_2 \geq \lambda_3 \geq \lambda_4 \geq \dots$$

- To compute the coefficients: $c_{ij} = (\vec{x}_i - \vec{\mu})^T \vec{v}_j$

Eigenfaces



Machine Learning 4771

Instructor: Tony Jebara

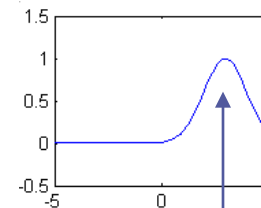
Topic 11

- Maximum Likelihood as Bayesian Inference
- Maximum A Posteriori
- Bayesian Gaussian Estimation

Why Maximum Likelihood?

- So far, assumed max (log) likelihood (IID or otherwise)

- Philosophical: Why? $\max_{\theta} L(\theta) = \max_{\theta} p(x_1, \dots, x_N | \theta)$
 $= \max_{\theta} \prod_{i=1}^N p(x_i | \theta)$



- Also, why ignore $p(\theta)$?

- Hint: Recall Bayes rule:

$$\begin{array}{c}
 \text{likelihood} \rightarrow \\
 \text{posterior} \rightarrow p(\theta | x) = \frac{p(x | \theta) p(\theta)}{p(x)} \leftarrow \text{prior} \\
 \hspace{15em} \leftarrow \text{evidence}
 \end{array}$$

- Everyone agrees on probability theory: inference and use of probability models when we have computed $p(x)$
- But how get to $p(x)$ from data? Debate...
- Two schools of thought: Bayesians and Frequentists

Bayesians & Frequentists

- Frequentists (Neymann/Pearson/Wald). An orthodox view that sampling is infinite and decision rules can be sharp.
- Bayesians (Bayes/Laplace/de Finetti). Unknown quantities are treated probabilistically and the state of the world can always be updated.



de Finetti: $p(\text{event}) = \text{price I would pay for a contract that pays 1\$ when event happens}$

- Likelihoodists (Fisher). Single sample inference based on maximizing the likelihood function and relying on the Birnbaum's Theorem. Bayesians – But they don't know it.

Bayesians & Frequentists

- Frequentists:
 - Data are a repeatable random sample- there is a frequency
 - Underlying parameters remain constant during this repeatable process
 - Parameters are fixed
- Bayesians:
 - Data are observed from the realized sample.
 - Parameters are unknown and described probabilistically
 - Data are fixed

Bayesians & Frequentists

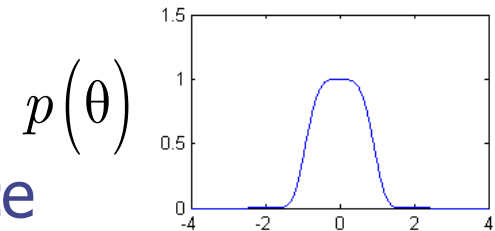
- **Frequentists:** classical / objective view / no priors
every statistician should compute same $p(x)$ so no priors
can't have a $p(\text{event})$ if it never happened
avoid $p(\theta)$, there is 1 true model, not distribution of them
permitted: $p_{\theta}(x,y)$ forbidden: $p(x,y|\theta)$
Frequentist inference: estimate one best model θ
use the **ML estimator** (unbiased & minimum variance)
do not depend on Bayes rule for learning
- **Bayesians:** subjective view / priors are ok
put a distribution or pdf on all variables in the problem
even models & deterministic quantities (i.e. speed of light)
use a prior $p(\theta)$, on the model θ before seeing any data
Bayesian inference: use Bayes rule for learning, integrate
over all model (θ) unknown variables

Bayesian Inference

- Bayes rule gives rise to maximum likelihood
- Assume we have a prior over models $p(\theta)$

$$\begin{array}{c}
 \text{likelihood} \rightarrow p(x | \theta) \\
 \text{posterior} \rightarrow p(\theta | x) = \frac{p(x | \theta) p(\theta)}{p(x)} \leftarrow \text{prior} \\
 \text{evidence} \leftarrow p(x)
 \end{array}$$

- How to pick $p(\theta)$?
 Pick simpler θ is better
 Pick form for mathematical convenience



- We have data (can assume IID): $\mathcal{X} = \{x_1, x_2, \dots, x_N\}$
- Want to get a model to compute: $p(x)$
- Want $p(x)$ given our data... How to proceed?

Bayesian Inference

- Want $p(x)$ given our data... $p(x | \mathcal{X}) = p(x | x_1, x_2, \dots, x_n)$

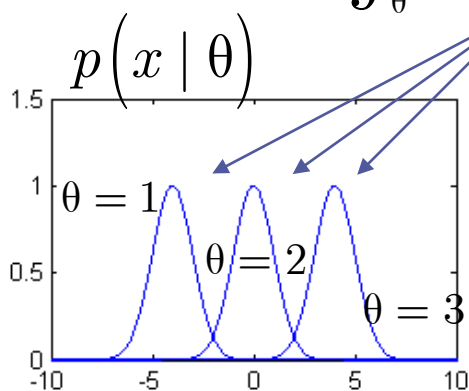
$$p(x | \mathcal{X}) = \int_{\theta} p(x, \theta | \mathcal{X}) d\theta$$

$$= \int_{\theta} p(x | \theta, \mathcal{X}) p(\theta | \mathcal{X}) d\theta$$

$$= \int_{\theta} p(x | \theta, \mathcal{X}) \frac{p(\mathcal{X} | \theta) p(\theta)}{p(\mathcal{X})} d\theta$$

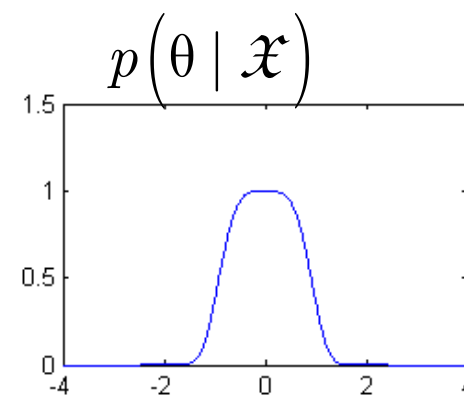
Prior \swarrow

$$= \int_{\theta} p(x | \theta) \frac{\prod_{i=1}^N p(x_i | \theta) p(\theta)}{p(\mathcal{X})} d\theta$$



**Many
models**

**Weight on
each model**



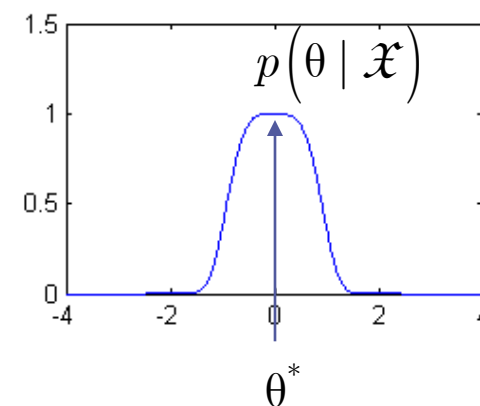
Bayesian Inference to MAP & ML

- The full **Bayesian Inference** integral can be mathematically tricky. Maximum likelihood is an approximation of it...

$$p(x | \mathcal{X}) = \int_{\theta} p(x | \theta) \frac{\prod_{i=1}^N p(x_i | \theta) p(\theta)}{p(\mathcal{X})} d\theta$$

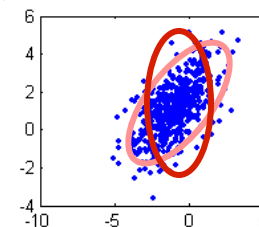
$$\approx \int_{\theta} p(x | \theta) \delta(\theta - \theta^*) d\theta$$

$$\text{where } \theta^* = \begin{cases} \arg \max_{\theta} \frac{\prod_{i=1}^N p(x_i | \theta) p(\theta)}{p(\mathcal{X})} & \text{MAP} \\ \arg \max_{\theta} \frac{\prod_{i=1}^N p(x_i | \theta) \text{uniform}(\theta)}{p(\mathcal{X})} & \text{ML} \end{cases}$$



- Maximum A Posteriori (MAP)** is like **Maximum Likelihood (ML)** with a prior $p(\theta)$ which lets us prefer some models over others

$$l_{MAP}(\theta) = l_{ML}(\theta) + \log p(\theta) = \sum_{i=1}^N \log p(x_i | \theta) + \log p(\theta)$$



Bayesian Inference Example

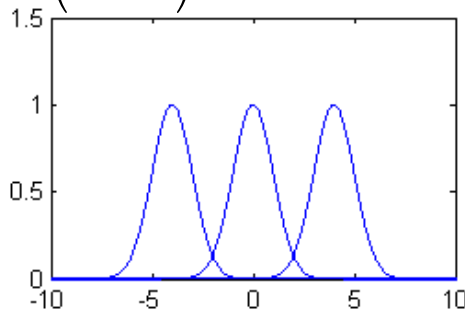
- For Gaussians, we CAN compute the integral (but hard!)

$$p(x | \mathcal{X}) = \int_{\theta} p(x | \theta) \frac{\prod_{i=1}^N p(x_i | \theta) p(\theta)}{p(\mathcal{X})} d\theta$$

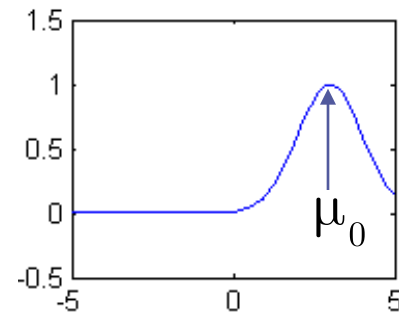
$$\propto \int_{\theta} p(x | \theta) \prod_{i=1}^N p(x_i | \theta) p(\theta) d\theta$$

- Example: ... assume 1d Gaussian & Gaussian prior on mean

$$p(x | \theta) = \text{Gaussian}$$



$$p(\theta) = \text{Gaussian}$$



$$p(x | \mathcal{X}) \propto \int_{\mu} \left(\frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x-\mu)^2} \right) \prod_{i=1}^N \left(\frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x_i-\mu)^2} \right) \left(\frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(\mu-\mu_0)^2} \right) d\mu$$

Bayesian Inference Example

- Solve integral over all Gaussian means with variance=1

$$\begin{aligned}
 p(x | \mathcal{X}) &\propto \int_{\mu=-\infty}^{\mu=\infty} \left(\frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x-\mu)^2} \right) \prod_{i=1}^N \left(\frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x_i-\mu)^2} \right) \left(\frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(\mu_0-\mu)^2} \right) d\mu \\
 &\propto \int_{\mu=-\infty}^{\mu=\infty} \exp \left(-\frac{1}{2}(x-\mu)^2 - \sum_i \frac{1}{2}(x_i-\mu)^2 - \frac{1}{2}(\mu_0-\mu)^2 \right) d\mu \\
 &\propto \int_{\mu=-\infty}^{\mu=\infty} \exp \left(-\frac{1}{2} \left[(N+2)\mu^2 - 2\mu \left(x + \mu_0 + \sum_i x_i \right) + x^2 \right] \right) d\mu \\
 &\propto \int_{\mu=-\infty}^{\mu=\infty} \exp \left(-\frac{1}{2} \left[(N+2)\mu^2 - 2\mu \left(x + \mu_0 + \sum_i x_i \right) + x^2 \right] + \left[\right]^2 - \left[\right]^2 \right) d\mu \\
 &\propto \exp \left(-\frac{1}{2} \left[\frac{-(x+\mu_0+\sum_i x_i)^2}{N+2} + x^2 \right] \right) \quad \tilde{\mu} = \frac{\mu_0 + \sum_i x_i}{N+1} \\
 &= N(x | \tilde{\mu}, \tilde{\sigma}^2) \quad \tilde{\sigma}^2 = \frac{N+2}{N+1}
 \end{aligned}$$

- Can integrate over μ and Σ for multivariate Gaussian (Jordan ch. 4 and Minka Tutorial)

$$p(x | \mathcal{X}) = \frac{\Gamma((N+1)/2)}{\Gamma((N+1-d)/2)} \left| \frac{1}{(N+1)\pi} \bar{\Sigma}^{-1} \right|^{1/2} \left(\frac{1}{N+1} (x - \bar{\mu})^T \bar{\Sigma}^{-1} (x - \bar{\mu}) + 1 \right)^{-(N+1)/2}$$

