

Piazza started

Announcements

Matlab Grader homework, email Friday,
2 (of 9) homeworks Due 21 April, Binary graded.

Jupyter homework?: translate matlab to Jupiter, TA Harshul h6gupta@eng.ucsd.edu or me
I would like this to happen.

“GPU” homework. NOAA climate data in Jupyter on the datahub.ucsd.edu, 15 April.

Projects: Any language

Podcast might work eventually.

Today:

- Stanford CNN
- Bernoulli
- Gaussian 1.2
- Gaussian 2.3
- Decision theory 1.5
- Information theory 1.6

Monday

Stanford CNN, Linear models for regression 3

Non-parametric method

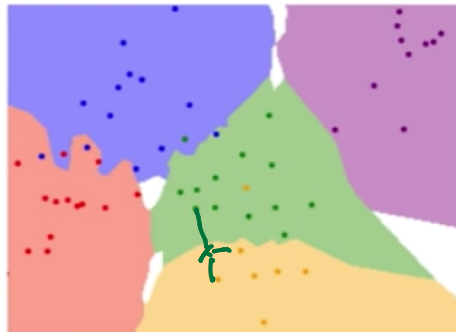
K-means

K-Nearest Neighbors

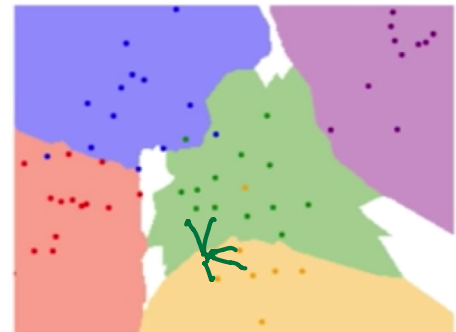
Instead of copying label from nearest neighbor,
take **majority vote** from K closest points



K = 1

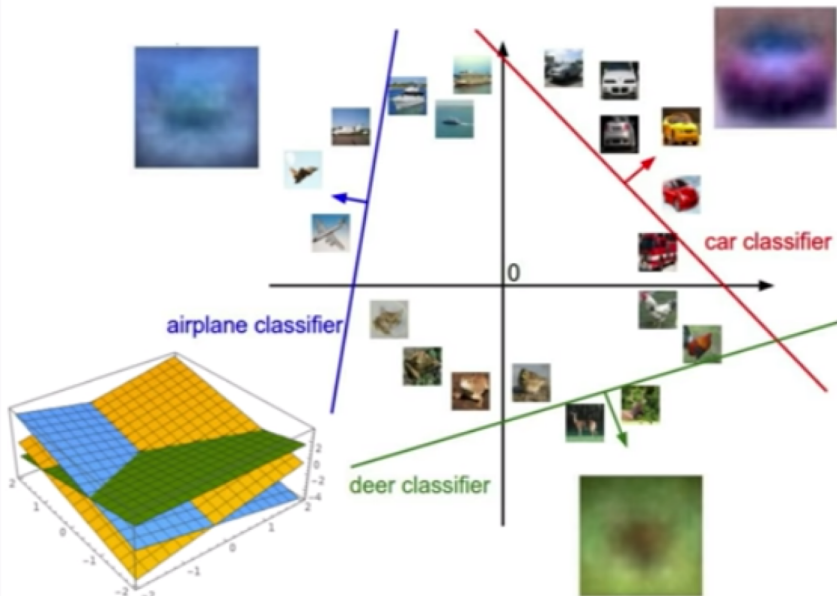


K = 3



K = 5

Interpreting a Linear Classifier



$$f(x, W) = Wx + b$$



Array of **32x32x3** numbers
(3072 numbers total)

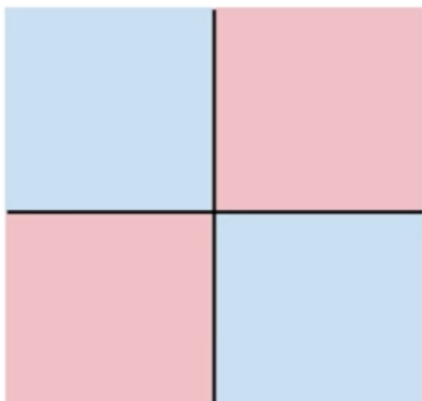
Hard cases for a linear classifier

Class 1:

number of pixels > 0 odd

Class 2:

number of pixels > 0 even

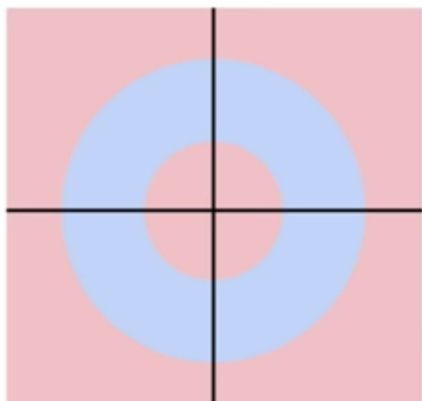


Class 1:

$1 \leq L2 \text{ norm} \leq 2$

Class 2:

Everything else

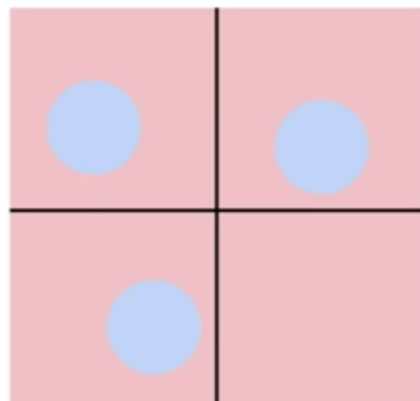


Class 1:

Three modes

Class 2:

Everything else



Coin estimate (Bishop 2.1)

- Binary variables $x \in \{0, 1\}$

$$p(x = 1 | \mu) = \mu$$

- Bernoulli distributed

$$\Rightarrow \text{Bern}(x | \mu) = \mu^x (1 - \mu)^{1-x}$$

$$\mathbb{E}[x] = \mu$$

$$\text{var}[x] = \mu(1 - \mu).$$

$$p(x=1) = p(H) = \mu$$

$$p(x=0) = p(T) = 1 - \mu$$

$$E(x) = \int x p(x) dx$$

$$(2.2) \quad = 1 \cdot \mu + 0 \cdot (1 - \mu) = \mu$$

$$\text{Var}(x) = E[(x - \bar{x})^2]$$

$$= (1 - \mu)^2 \cdot \mu + (0 - \mu)^2 (1 - \mu)$$

- N observations, Likelihood:

$$p(\mathcal{D} | \mu) = \prod_{n=1}^N p(x_n | \mu) = \prod_{n=1}^N \mu^{x_n} (1 - \mu)^{1-x_n}. \quad (2.5)$$

$$\ell = \ln p(\mathcal{D} | \mu) = \sum_{n=1}^N \ln p(x_n | \mu) = \sum_{n=1}^N \{x_n \ln \mu + (1 - x_n) \ln(1 - \mu)\}. \quad (2.6)$$

- Max likelihood

$$\frac{\partial \ell}{\partial \mu} = \sum \frac{x_n}{\mu} + \frac{1 - x_n}{1 - \mu}$$

$$\mu_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N x_n$$

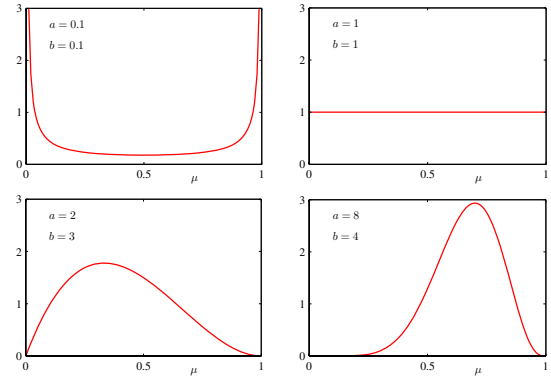
Coin estimate (Bishop 2.1)

post. like prior

- Bayes $p(x|y)=p(y|x)p(x)$

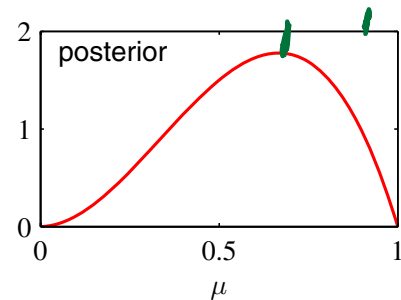
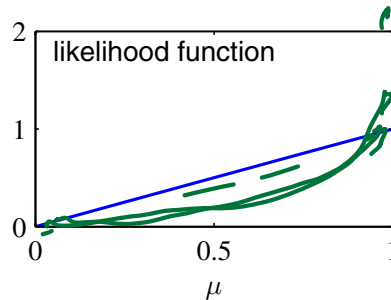
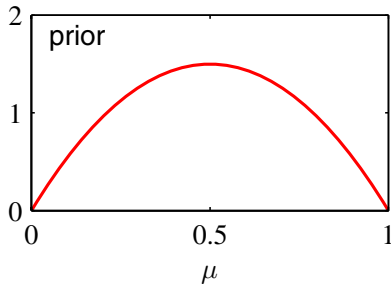
- Conjugate prior

$$\text{Beta}(\mu|a, b) = \frac{\Gamma(a + b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1 - \mu)^{b-1}$$



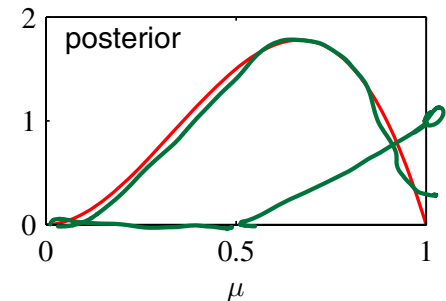
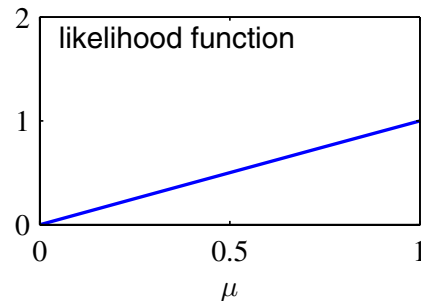
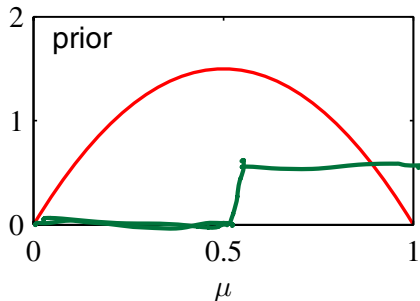
Bayes:

b = a = 2



ML MAP BAYES

- ML point estimate
- MAP point estimate (often in literature ML=MAP)
- Bayes => probability => From which all information can be obtained
 - MAP, median, error estimates
 - Further analysis as sequential →.
 - Disadvantage... not a point estimate.



Bayes Rule

$$P(\text{hypothesis}|\text{data}) = \frac{P(\text{data}|\text{hypothesis})P(\text{hypothesis})}{P(\text{data})}$$

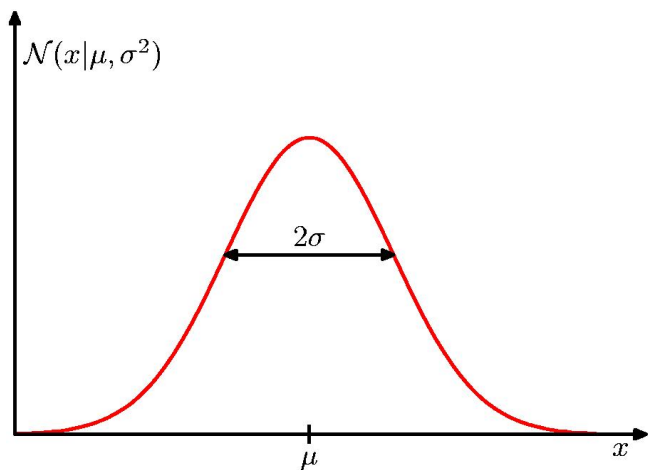


Rev'd Thomas Bayes (1702–1761)

- Bayes rule tells us how to do inference about hypotheses from data.
- Learning and prediction can be seen as forms of inference.

The Gaussian Distribution

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left\{-\frac{1}{2\sigma^2}(x - \mu)^2\right\}$$



$$\underline{\mathcal{N}(x|\mu, \sigma^2) > 0}$$

$$\underline{\int_{-\infty}^{\infty} \mathcal{N}(x|\mu, \sigma^2) dx = 1}$$

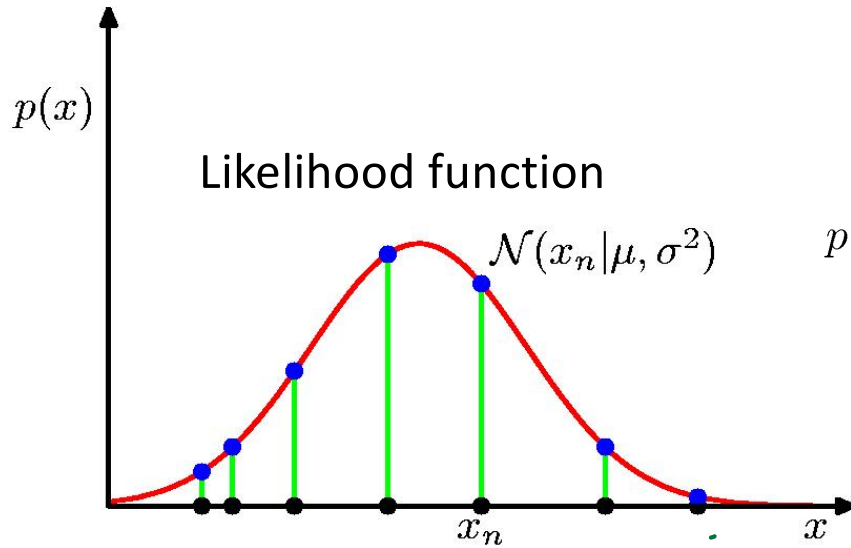
Gaussian Mean and Variance

$$\mathbb{E}[x] = \int_{-\infty}^{\infty} \mathcal{N}(x|\mu, \sigma^2) x dx = \underline{\mu}$$

$$\mathbb{E}[x^2] = \int_{-\infty}^{\infty} \mathcal{N}(x|\mu, \sigma^2) x^2 dx = \mu^2 + \sigma^2$$

$$\text{var}[x] = \mathbb{E}[x^2] - \mathbb{E}[x]^2 = \underline{\sigma^2}$$

Gaussian Parameter Estimation



$$p(\mathbf{x} | \mu, \sigma^2) = \prod_{n=1}^N \mathcal{N}(x_n | \mu, \sigma^2)$$

$$\frac{\partial \ell}{\partial \mu} = \frac{1}{2\sigma^2} \sum_n^N (x_n - \mu) = 0$$

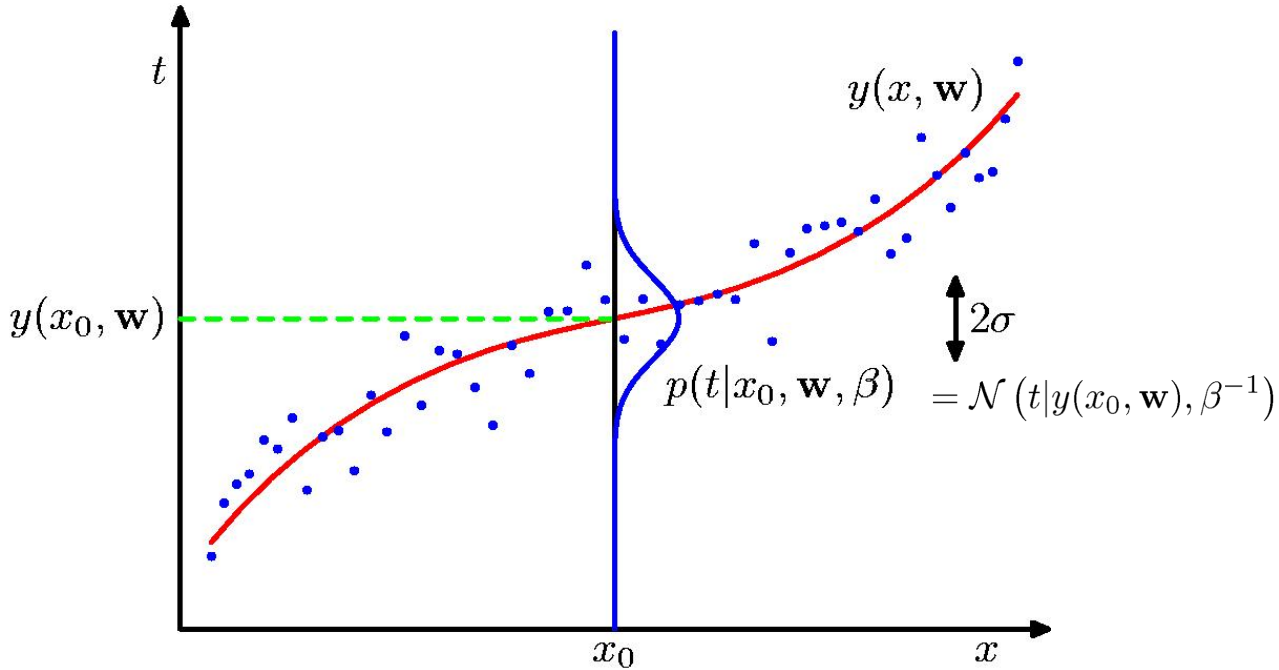
Maximum (Log) Likelihood

$$\ln p(\mathbf{x} | \mu, \sigma^2) = -\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2 - \frac{N}{2} \ln \sigma^2 - \frac{N}{2} \ln(2\pi)$$

$$\mu_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N x_n$$

$$\sigma_{\text{ML}}^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{\text{ML}})^2$$

Curve Fitting Re-visited, Bishop1.2.5



Maximum Likelihood

$$p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n|y(x_n, \mathbf{w}), \beta^{-1}). \quad (1.61)$$

As we did in the case of the simple Gaussian distribution earlier, it is convenient to maximize the logarithm of the likelihood function. Substituting for the form of the Gaussian distribution, given by (1.46), we obtain the log likelihood function in the form

$$\ln p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) = -\frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi). \quad (1.62)$$

$$\frac{1}{\beta_{\text{ML}}} = \frac{1}{N} \sum_{n=1}^N \{y(x_n, \mathbf{w}_{\text{ML}}) - t_n\}^2. \quad (1.63)$$

Giving estimates of \mathbf{w} and β , we can predict

$$\underline{p(t|x, \mathbf{w}_{\text{ML}}, \beta_{\text{ML}})} = \mathcal{N}(t|y(x, \mathbf{w}_{\text{ML}}), \beta_{\text{ML}}^{-1}). \quad (1.64)$$

MAP: A Step towards Bayes 1.2.5

prior

$$p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I}) = \underbrace{\left(\frac{\alpha}{2\pi}\right)^{(M+1)/2}}_{\text{prior}} \exp\left\{-\frac{\alpha}{2}\mathbf{w}^T\mathbf{w}\right\}$$

$$p(\mathbf{w}|\mathbf{x}, \mathbf{t}, \alpha, \beta) \propto p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta)p(\mathbf{w}|\alpha)$$

$$\mathcal{N}(\mathbf{w}^T\mathbf{x}, \beta^{-1})\mathcal{N}(\mathbf{0}, \alpha)$$

$$-\ln(p(\mathbf{w}|\mathbf{t})) = \frac{\beta}{2} \sum_n^N (\mathbf{w}^T\mathbf{x}_n - t_n)^2 + \frac{\alpha}{2}\mathbf{w}^T\mathbf{w} + \text{const}$$

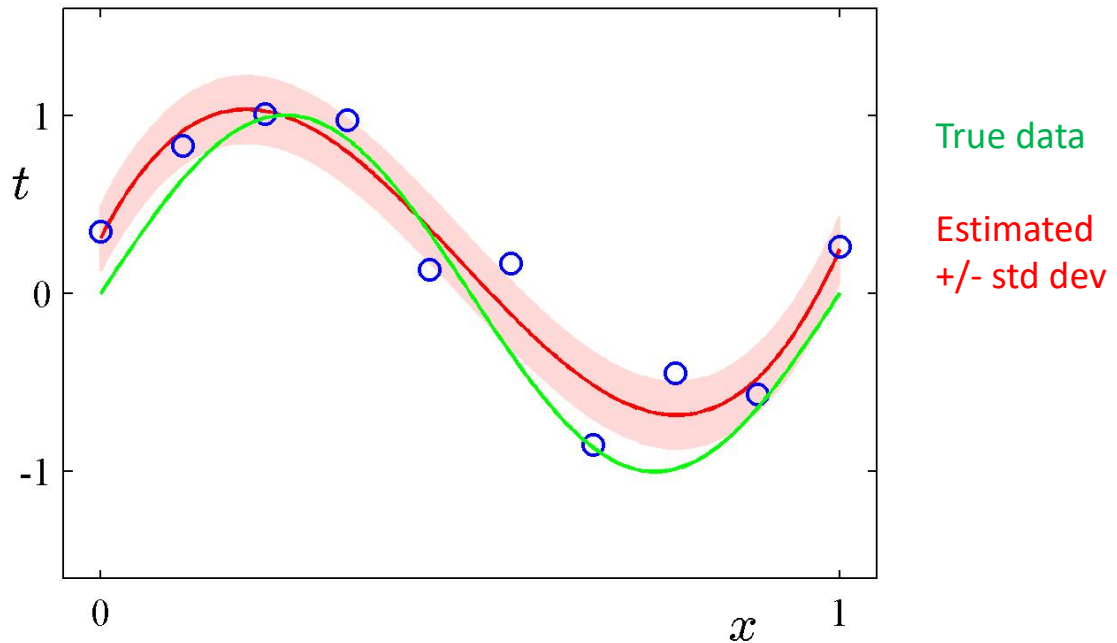
$$\rightarrow \beta\tilde{E}(\mathbf{w}) = \frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{\alpha}{2}\underbrace{\mathbf{w}^T\mathbf{w}}_{\|\mathbf{w}\|_2^2}$$

Determine \mathbf{w}_{MAP} by minimizing regularized sum-of-squares error, $\tilde{E}(\mathbf{w})$.

Regularized sum of squares \leftarrow

Predictive Distribution

$$p(t|x, \mathbf{w}_{\text{ML}}, \beta_{\text{ML}}) = \mathcal{N}(t|y(x, \mathbf{w}_{\text{ML}}), \beta_{\text{ML}}^{-1})$$



Parametric Distributions

Basic building blocks:

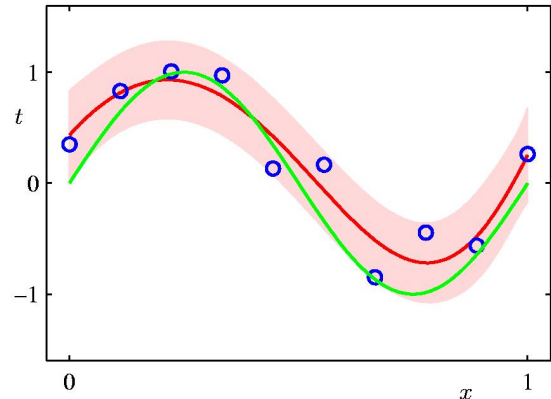
Need to determine θ given $p(\mathbf{x}|\theta)$
 $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$

Representation: θ^* or $p(\theta)$

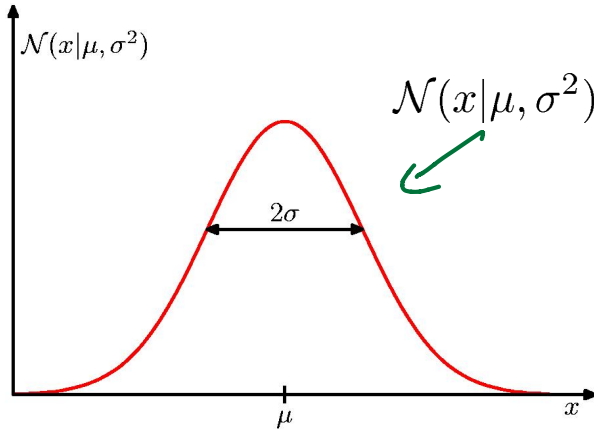
Recall Curve Fitting

$$p(t|x, \mathbf{x}, \mathbf{t}) = \int p(t|x, \mathbf{w})p(\mathbf{w}|\mathbf{x}, \mathbf{t}) d\mathbf{w}$$

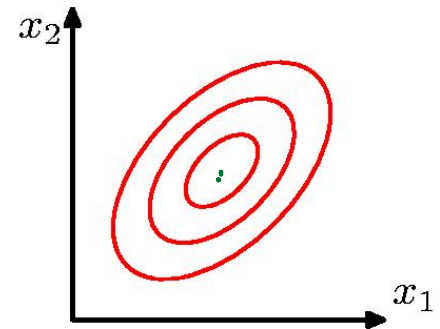
We focus on Gaussians!



The Gaussian Distribution



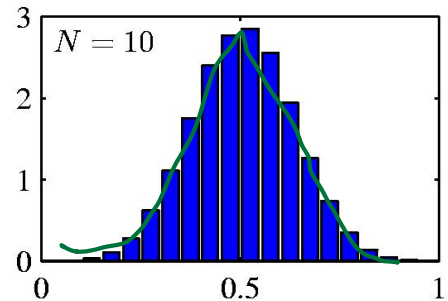
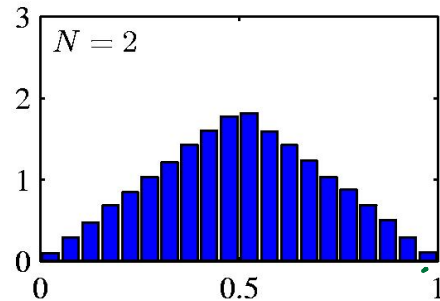
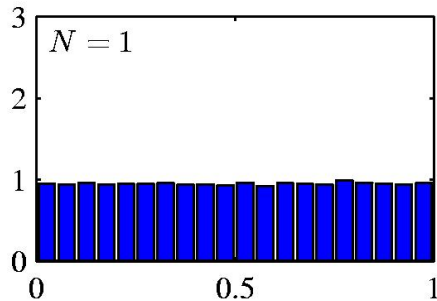
$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left\{-\frac{1}{2\sigma^2}(x - \mu)^2\right\}$$



$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right\}$$

Central Limit Theorem

- The distribution of the sum of N i.i.d. random variables becomes increasingly Gaussian as N grows.
- Example: N uniform $[0,1]$ random variables.



Geometry of the Multivariate Gaussian

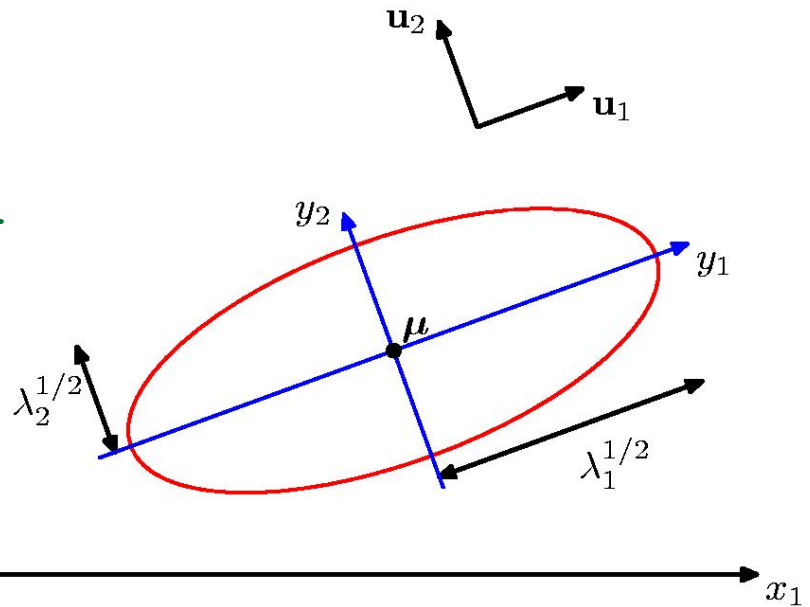
$$e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})}$$

$$\Delta^2 = (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})$$

$$\boldsymbol{\Sigma}^{-1} = \sum_{i=1}^D \frac{1}{\lambda_i} \mathbf{u}_i \mathbf{u}_i^T$$

$$\Delta^2 = \sum_{i=1}^D \frac{y_i^2}{\lambda_i} = \mathbf{y}^T \boldsymbol{\Lambda}^{-1} \mathbf{y}$$

$$y_i = \mathbf{u}_i^T (\mathbf{x} - \boldsymbol{\mu})$$



Moments of the Multivariate Gaussian (2)

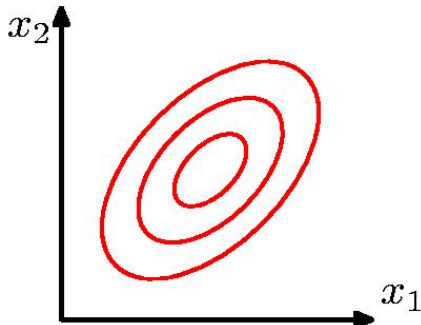
$$\mathbb{E}[\mathbf{x}] = \boldsymbol{\mu}$$

$$\mathbb{E}[\mathbf{x}\mathbf{x}^T] = \boldsymbol{\mu}\boldsymbol{\mu}^T + \boldsymbol{\Sigma}$$

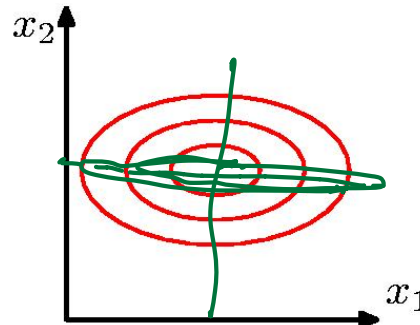
$$\text{cov}[\mathbf{x}] = \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^T] = \boldsymbol{\Sigma} = \begin{bmatrix} \sigma_1^2 & & \\ & \sigma_2^2 & \\ & & \sigma_N^2 \end{bmatrix}$$

A Gaussian requires $\underline{D*(D-1)/2 + D}$ parameters.
Often we use $\underline{D + D}$ or
Just $D+1$ parameters.

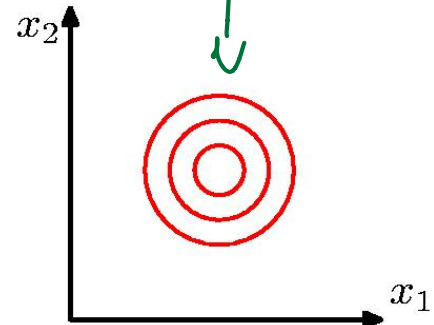
$$= \sigma^2 \mathbf{I}$$



(a)



(b)



(c)

Partitioned Conditionals and Marginals, page 89

Conditionals

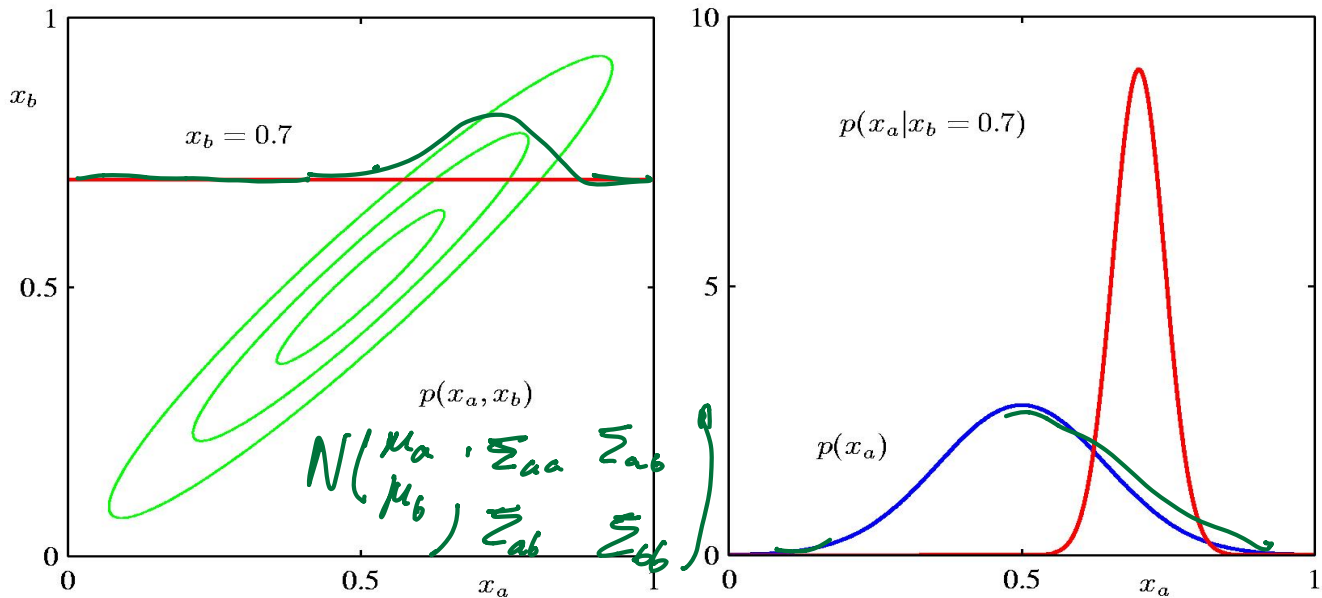
$$p(\mathbf{x}_a | \mathbf{x}_b) = \mathcal{N}(\mathbf{x}_a | \boldsymbol{\mu}_{a|b}, \boldsymbol{\Sigma}_{a|b})$$

$$\boldsymbol{\Sigma}_{a|b} = \boldsymbol{\Lambda}_{aa}^{-1} = \boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab} \boldsymbol{\Sigma}_{bb}^{-1} \boldsymbol{\Sigma}_{ba}$$

$$\begin{aligned} \boldsymbol{\mu}_{a|b} &= \boldsymbol{\Sigma}_{a|b} \{ \boldsymbol{\Lambda}_{aa} \boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b) \} \\ &= \boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{aa}^{-1} \boldsymbol{\Lambda}_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b) \\ &= \boldsymbol{\mu}_a + \boldsymbol{\Sigma}_{ab} \boldsymbol{\Sigma}_{bb}^{-1} (\mathbf{x}_b - \boldsymbol{\mu}_b) \end{aligned}$$

marginal

$$\begin{aligned} \underline{p(\mathbf{x}_a)} &= \int p(\mathbf{x}_a, \mathbf{x}_b) d\mathbf{x}_b \\ &= \underline{\mathcal{N}(\mathbf{x}_a | \boldsymbol{\mu}_a, \boldsymbol{\Sigma}_{aa})} \end{aligned}$$



ML for the Gaussian (1) Bishop 2.3.4

Given i.i.d. data $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)^T$, the log likelihood function is given by

$$\ln p(\mathbf{X} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln |\boldsymbol{\Sigma}| - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu})$$

$$\hookrightarrow \ell = -\ln p$$

$$= \frac{N}{2} \ln |\boldsymbol{\Sigma}| + \frac{1}{2} \text{Tr} \left(\sum_n (\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) \right)$$

$$= \frac{N}{2} \ln |\boldsymbol{\Sigma}| + \frac{1}{2} \text{tr} \left(\sum_n (\mathbf{x}_n - \boldsymbol{\mu}) (\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} \right)$$

$$= \frac{N}{2} \left[\ln |\boldsymbol{\Sigma}| + \text{tr} (\mathbf{S}_y \boldsymbol{\Sigma}^{-1}) \right]$$

$$\mathbf{S}_y = \frac{1}{N} \sum (\mathbf{x} - \boldsymbol{\mu}) (\mathbf{x} - \boldsymbol{\mu})^T$$

$$\frac{\partial \ell}{\partial \boldsymbol{\mu}}$$

$$\frac{\partial \ell}{\partial \boldsymbol{\mu}} = \boldsymbol{\Sigma}^{-1} + \mathbf{S}_y \boldsymbol{\Sigma}^{-1} \boldsymbol{\Sigma}^{-1} = 0$$

$$\boldsymbol{\Sigma} = \mathbf{S}_\mu$$

$$\text{tr}(\mathbf{A}\mathbf{B}\mathbf{C}) = \text{tr}(\mathbf{C}\mathbf{A}\mathbf{B})$$

$$\frac{\partial}{\partial \mathbf{A}} \ln |\mathbf{A}| = (\mathbf{A}^{-1})^T \quad (\text{C.28})$$

$$\frac{\partial}{\partial \mathbf{A}} \text{Tr}(\mathbf{A}\mathbf{B}) = \mathbf{B}^T. \quad (\text{C.24})$$

$$\frac{\partial}{\partial x} (\mathbf{A}^{-1}) = -\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x} \mathbf{A}^{-1} \quad (\text{C.21})$$

Maximum Likelihood for the Gaussian

- Set the derivative of the log likelihood function to zero,

- and solve to obtain
$$\frac{\partial}{\partial \boldsymbol{\mu}} \ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \boldsymbol{\Sigma}^{-1}(\mathbf{x}_n - \boldsymbol{\mu}) = 0$$

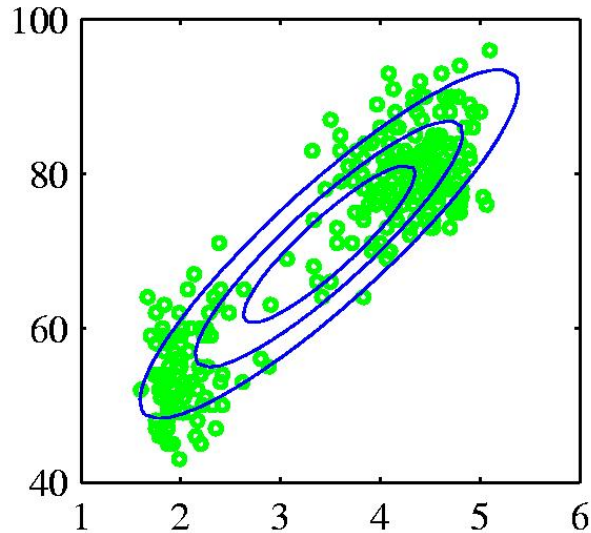
- Similarly
$$\boldsymbol{\mu}_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n.$$

$$\boldsymbol{\Sigma}_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu}_{\text{ML}})(\mathbf{x}_n - \boldsymbol{\mu}_{\text{ML}})^{\text{T}}.$$

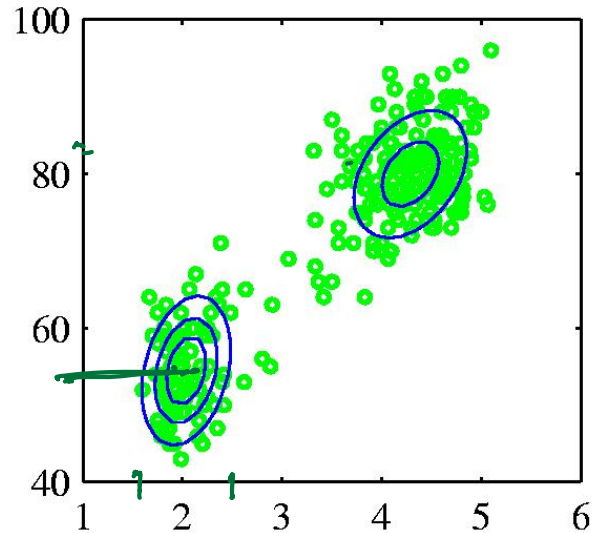
Mixtures of Gaussians (Bishop 2.3.9)

Old Faithful geyser:

The time between eruptions has a [bimodal distribution](#), with the mean interval being either 65 or 91 minutes, and is dependent on the length of the prior eruption. Within a margin of error of ± 10 minutes, Old Faithful will erupt either 65 minutes after an eruption lasting less than $2\frac{1}{2}$ minutes, or 91 minutes after an eruption lasting more than $2\frac{1}{2}$ minutes.



Single Gaussian



Mixture of two Gaussians

Mixtures of Gaussians (Bishop 2.3.9)

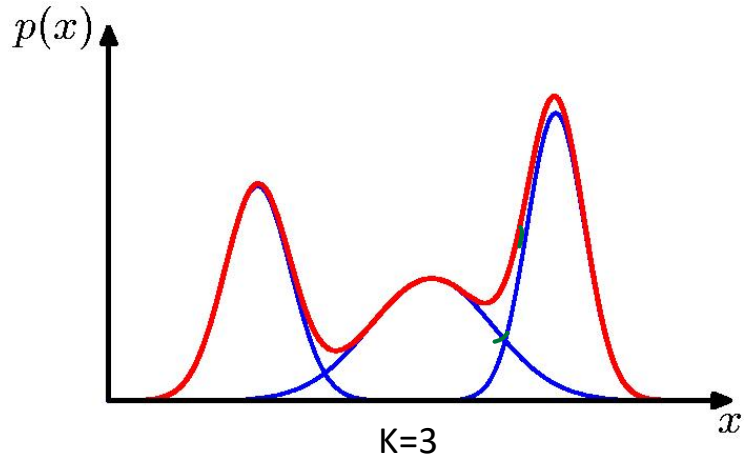
- Combine simple models into a complex model:

$$\underline{p(\mathbf{x})} = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

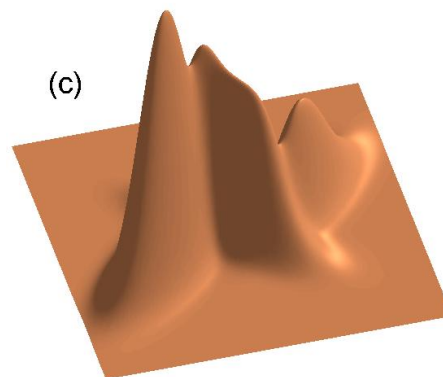
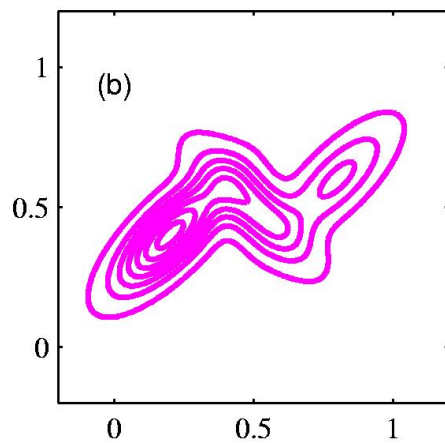
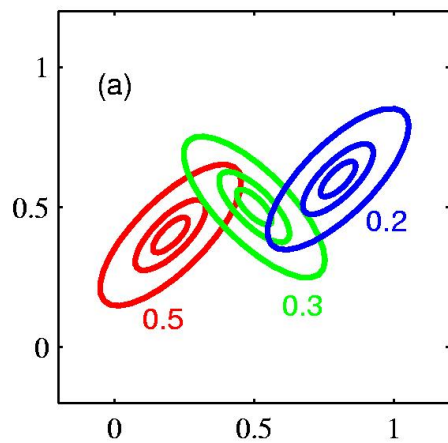
↑
Mixing coefficient

Component

$$\underline{\forall k : \pi_k \geq 0} \quad \underline{\sum_{k=1}^K \pi_k = 1}$$




Mixtures of Gaussians (Bishop 2.3.9)



Mixtures of Gaussians (Bishop 2.3.9)

- Determining parameters $\boldsymbol{\pi}$, $\boldsymbol{\mu}$, and $\boldsymbol{\Sigma}$ using maximum log likelihood

$$\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$


Log of a sum; no closed form maximum.

- Solution: use standard, iterative, numeric optimization methods or the *expectation maximization* algorithm (Chapter 9).

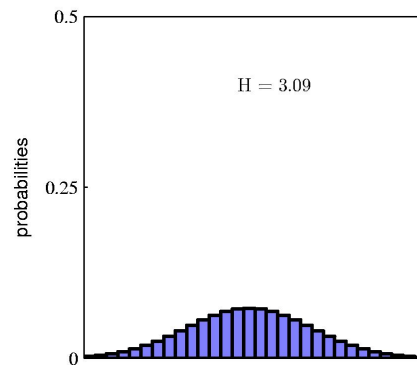
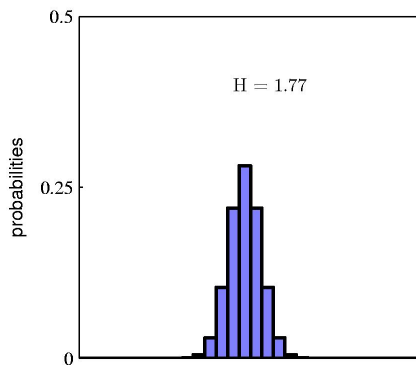
EM

Entropy 1.6

$$H[x] = - \sum_x p(x) \log_2 p(x)$$

Important quantity in

- coding theory
- statistical physics
- machine learning



Differential Entropy

Put bins of width Δ along the real line

$$\lim_{\Delta \rightarrow 0} \left\{ - \sum_i p(x_i) \Delta \ln p(x_i) \right\} = - \int p(x) \ln p(x) dx$$

For fixed σ^2 differential entropy maximized when

in which case

$$p(x) = \mathcal{N}(x|\mu, \sigma^2)$$

$$H[x] = \frac{1}{2} \{1 + \ln(2\pi\sigma^2)\}.$$

The Kullback-Leibler Divergence

P true distribution, q is approximating distribution

$$\begin{aligned}\text{KL}(p\|q) &= - \int p(\mathbf{x}) \ln q(\mathbf{x}) \, d\mathbf{x} - \left(- \int p(\mathbf{x}) \ln p(\mathbf{x}) \, d\mathbf{x} \right) \\ &= - \int p(\mathbf{x}) \ln \left\{ \frac{q(\mathbf{x})}{p(\mathbf{x})} \right\} \, d\mathbf{x}\end{aligned}$$

$$\text{KL}(p\|q) \simeq \frac{1}{N} \sum_{n=1}^N \{-\ln q(\mathbf{x}_n|\boldsymbol{\theta}) + \ln p(\mathbf{x}_n)\}$$

$$\text{KL}(p\|q) \geq 0$$

$$\text{KL}(p\|q) \neq \text{KL}(q\|p)$$