

Announcements



✓ **Matlab Grader homework**, emailed Thursday,
1 (of 9) homeworks Due 21 April, Binary graded.
2 this week

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 computer language

Podcast might work eventually.

Today:

- Stanford CNN
- Gaussian, Bishop 2.3
- Gaussian Process 6.4
- Linear regression 3.0-3.2

Wednesday 10 April

Stanford CNN, Linear models for regression 3, Applications of Gaussian processes.

Bayes and Softmax (Bishop p. 198)

- Bayes:

$$p(x|y) = \frac{p(y|x)p(x)}{p(y)} = \frac{p(y|x)p(x)}{\sum_{y \in Y} p(x, y)}$$

- Classification of N classes:

$$p(\mathcal{C}_n | \mathbf{x}) = \frac{p(\mathbf{x} | \mathcal{C}_n) p(\mathcal{C}_n)}{\sum_{k=1}^N p(\mathbf{x} | \mathcal{C}_k) p(\mathcal{C}_k)}$$

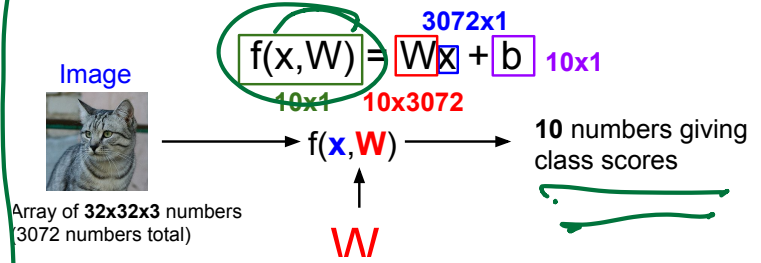
$$= \frac{\exp(a_n)}{\sum_{k=1}^N \exp(a_k)}$$

$\sum_{k=1}^N$
 $\sum_{k=1}^N$

with

$$a_n = \ln(p(\mathbf{x} | \mathcal{C}_n) p(\mathcal{C}_n))$$

Parametric Approach: Linear Classifier



Softmax to Logistic Regression (Bishop p. 198)

$$\begin{aligned} p(\mathcal{C}_1|\mathbf{x}) &= \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{\sum_{k=1}^2 p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)} \\ &= \frac{\exp(a_1)}{\sum_{k=1}^2 \exp(a_k)} = \frac{1}{1 + \underline{\exp(-a)}} \end{aligned}$$

with

$$a = \ln \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)}$$

Softmax with Gaussian (Bishop p. 198)

$$p(C_n | \mathbf{x}) = \frac{p(\mathbf{x} | C_n) p(C_n)}{\sum_{k=1}^N p(\mathbf{x} | C_k) p(C_k)}$$

$$= \frac{\exp(a_n)}{\sum_{k=1}^N \exp(a_k)}$$

with

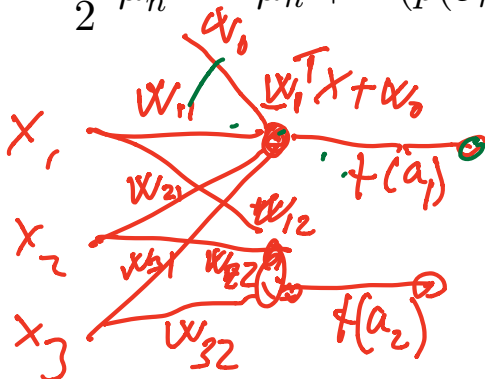
$$a_n = \ln(p(\mathbf{x} | C_n) p(C_n))$$

Assuming \mathbf{x} is Gaussian $\mathcal{N}(\mu_n, \Sigma)$

$$a_n = \mathbf{w}_n^T \mathbf{x} + w_0$$

$$\mathbf{w}_n = \Sigma^{-1} \mu_n$$

$$w_0 = \frac{-1}{2} \mu_n^T \Sigma^{-1} \mu_n + \ln(p(C_n))$$



$$\mathbf{x} \in \mathcal{N}(\mu_n, \Sigma)$$

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

$$-\ln(p(\mathbf{x} | C_n)) \sim \frac{1}{2} \mathbf{x}^T \Sigma^{-1} \mathbf{x} + \frac{1}{2} \mu_n^T \Sigma^{-1} \mu_n - \mu_n^T \Sigma^{-1} \mathbf{x}$$

$\underbrace{\quad}_{\mathbf{w}_n}$

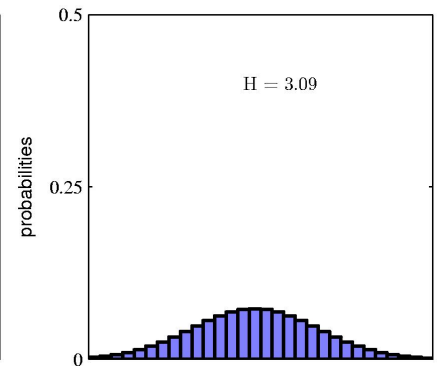
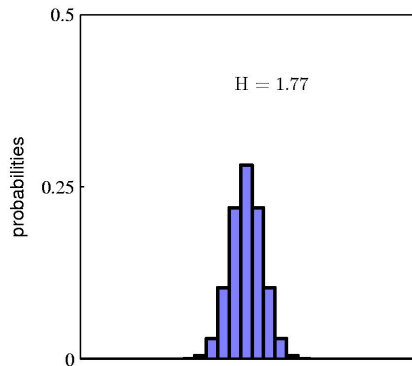
$$w_0 = \frac{1}{2} \mu_n^T \Sigma^{-1} \mu_n + \ln C_n$$

Entropy 1.6

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

Important quantity in

- coding theory
- statistical physics
- machine learning



The Kullback-Leibler Divergence

P true distribution, q is approximating distribution

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

not a distance measure

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

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

KL homework

- Support of P and Q = > “only >0” don't use isnan isinf
- After you pass. Take your time to clean up. Get close to 50

Lecture 3

- Homework
- Pod-cast lecture on-line
- Next lectures:
 - I posted a rough plan.
 - It is flexible though so please come with suggestions

Bayes for linear model

$$\underline{y} = \underline{Ax} + \underline{n} \quad n \sim N(\mathbf{0}, \mathbf{C}_n) \quad \underline{y} \sim N(\underline{Ax}, \mathbf{C}_n) \quad \text{prior: } \underline{x} \sim N(\mathbf{0}, \mathbf{C}_x)$$

$$p(\mathbf{x}|\mathbf{y}) \sim p(\mathbf{y}|\mathbf{x})p(\mathbf{x}) \sim N(\mathbf{x}_p, \mathbf{C}_p)$$

mean $\mathbf{x}_p = \mathbf{C}_p \mathbf{A}^T \mathbf{C}_n^{-1} \mathbf{y}$

Covariance $\mathbf{C}_p^{-1} = \mathbf{A}^T \mathbf{C}_n^{-1} \mathbf{A} + \mathbf{C}_x^{-1}$

$$\propto e^{-\frac{1}{2}(\mathbf{x} - \mathbf{x}_p)^T \mathbf{C}_p^{-1} (\mathbf{x} - \mathbf{x}_p)}$$

$$\propto e^{-\frac{1}{2}(\mathbf{Ax} - \mathbf{y})^T \mathbf{C}_n^{-1} (\mathbf{Ax} - \mathbf{y})} e^{-\frac{1}{2}\mathbf{x}^T \mathbf{C}_x^{-1} \mathbf{x}}$$

$$\propto e^{-\frac{1}{2} \underbrace{[\mathbf{x}^T \mathbf{A}^T \mathbf{C}_n^{-1} \mathbf{A} \mathbf{x} + \mathbf{x}^T \mathbf{C}_x^{-1} \mathbf{x}]_{\mathbf{x}^T \mathbf{C}_p^{-1} \mathbf{x}} - \mathbf{x}^T \mathbf{A}^T \mathbf{C}_n^{-1} \mathbf{y}}_{\mathbf{x}^T \mathbf{C}_p^{-1} \mathbf{x}}$$

Bayes' Theorem for Gaussian Variables

- Given

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \underline{\boldsymbol{\mu}}, \boldsymbol{\Lambda}^{-1})$$

- we have

$$p(\mathbf{y} | \mathbf{x}) = \mathcal{N}(\mathbf{y} | \mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1})$$

$$\rightarrow p(\mathbf{y}) = \mathcal{N}(\mathbf{y} | \mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^T)$$

$$\rightarrow p(\mathbf{x} | \mathbf{y}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\Sigma}\{\mathbf{A}^T\mathbf{L}(\mathbf{y} - \mathbf{b}) + \boldsymbol{\Lambda}\boldsymbol{\mu}\}, \boldsymbol{\Sigma})$$

- where

$$\boldsymbol{\Sigma} = (\boldsymbol{\Lambda} + \mathbf{A}^T\mathbf{L}\mathbf{A})^{-1}$$

Sequential Estimation of mean (Bishop 2.3.5)

Contribution of the N^{th} data point, \mathbf{x}_N

$$\begin{aligned}\underline{\mu}_{\text{ML}}^{(N)} &= \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \\ &= \frac{1}{N} \mathbf{x}_N + \frac{1}{N} \sum_{n=1}^{N-1} \mathbf{x}_n \\ &= \frac{1}{N} \mathbf{x}_N + \frac{N-1}{N} \mu_{\text{ML}}^{(N-1)} \\ &= \underline{\mu}_{\text{ML}}^{(N-1)} + \frac{1}{N} (\mathbf{x}_N - \mu_{\text{ML}}^{(N-1)})\end{aligned}$$

correction given \mathbf{x}_N
correction weight
old estimate

Bayesian Inference for the Gaussian (Bishop2.3.6)

Assume σ^2 is known. Given i.i.d. data

the likelihood function for μ is given by $\mathbf{x} = \{x_1, \dots, x_N\}$

$$L(\mu) \quad p(\mathbf{x}|\mu) = \prod_{n=1}^N p(x_n|\mu) = \frac{1}{(2\pi\sigma^2)^{N/2}} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2 \right\}.$$

- This has a Gaussian shape as a function of μ (but it is *not* a distribution over μ).

Bayesian Inference for the Gaussian (Bishop2.3.6)

- Combined with a **Gaussian prior over μ** , $p(\mu) = \mathcal{N}(\mu | \underline{\mu_0}, \underline{\sigma_0^2})$.
- this gives the posterior

$$p(\mu | \mathbf{x}) = \mathcal{N}(\mu | \underline{\mu_N}, \underline{\sigma_N^2})$$

$$p(\mu | \mathbf{x}) \propto p(\mathbf{x} | \mu) p(\mu).$$

$$\underline{\mu_N} = \frac{\sigma^2}{N\sigma_0^2 + \sigma^2} \mu_0 + \frac{N\sigma_0^2}{N\sigma_0^2 + \sigma^2} \mu_{\text{ML}}, \quad \mu_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N x_n$$

$$\underline{\frac{1}{\sigma_N^2}} = \underline{\frac{1}{\sigma_0^2} + \frac{N}{\sigma^2}}$$

$$\frac{(\mu - \mu_N)^2}{2\sigma_N^2} = \frac{1}{2\sigma^2} (x - \mu)^2 + \frac{(\mu - \mu_0)^2}{2\sigma_0^2}$$

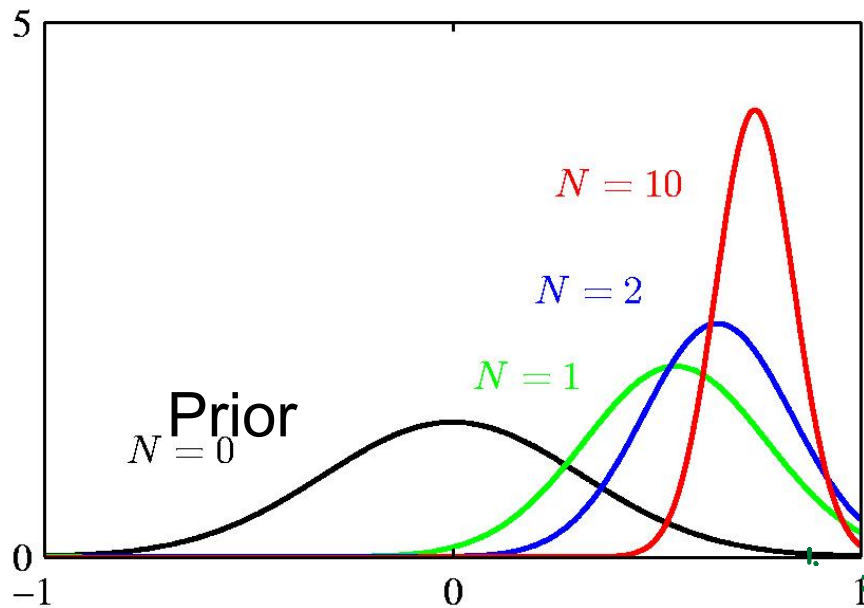
	<u>$N = 0$</u>	$N \rightarrow \infty$
μ_N	μ_0	μ_{ML}
σ_N^2	σ_0^2	0



Bayesian Inference for the Gaussian (3)

- Example: for $N = 0, 1, 2$ and 10 .

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



Bayesian Inference for the Gaussian (4)

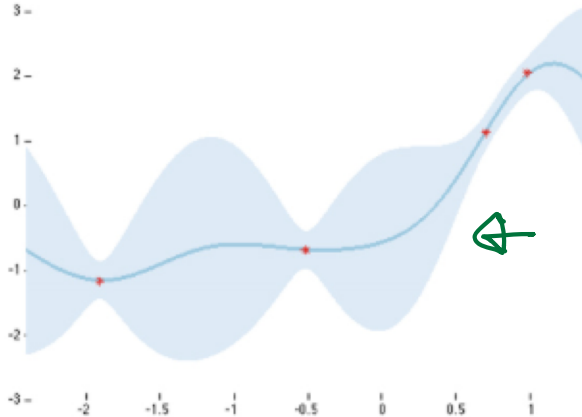
Sequential Estimation

$$\begin{aligned} \underline{p(\mu|\mathbf{x})} &\propto p(\mu)p(\mathbf{x}|\mu) \\ &= \left[p(\mu) \prod_{n=1}^{N-1} p(x_n|\mu) \right] p(x_N|\mu) \\ &\propto \mathcal{N}(\mu|\mu_{N-1}, \sigma_{N-1}^2) p(x_N|\mu) \end{aligned}$$

The posterior obtained after observing N-1 data points becomes the prior when we observe the Nth data point.

Conjugate prior: posterior and prior are in the same family. The **prior** is called a **conjugate prior** for the likelihood function.

Gaussian Process (Bishop 6.4, Murphy15)



$$t_n = y_n + \epsilon_n$$

$$f(\mathbf{x}) \sim GP(m(\mathbf{x}), \kappa(\mathbf{x}, \mathbf{x}'))$$

This is what a Gaussian process *posterior* looks like with 4 data points and a squared exponential covariance function. The bold blue line is the predictive mean, while the light blue shade is the predictive uncertainty (2 standard deviations). The model uncertainty is small near the data, and increases as we move away from the data points.

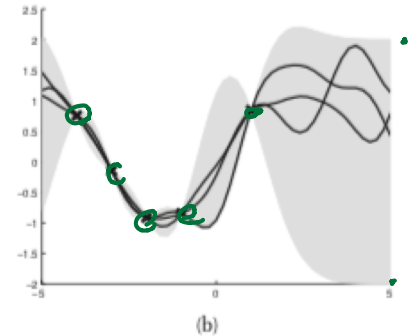
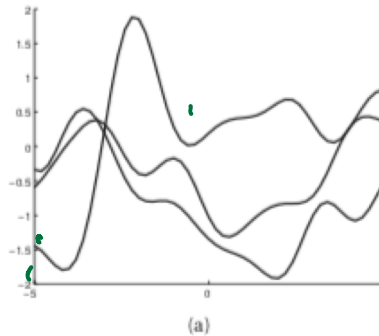


Figure 15.2 Left: some functions sampled from a GP prior with SE kernel. Right: some samples from a GP posterior, after conditioning on 5 noise-free observations. The shaded area represents $\mathbb{E}[f(\mathbf{x})] \pm 2\text{std}(f(\mathbf{x}))$. Based on Figure 2.2 of (Rasmussen and Williams 2006). Figure generated by `gprDemoNoiseFree`.

Gaussian Process (Murphy ch15)

$$f(\mathbf{x}) \sim \underline{GP}(m(\mathbf{x}), \kappa(\mathbf{x}, \mathbf{x}')) \quad \begin{aligned} \underline{m}(\mathbf{x}) &= \mathbb{E}[f(\mathbf{x})] \\ \underline{\kappa}(\mathbf{x}, \mathbf{x}') &= \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))^T] \end{aligned}$$

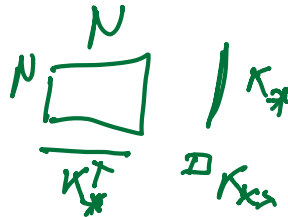
$$p(\mathbf{f}|\mathbf{X}) = \mathcal{N}(\mathbf{f}|\underline{\mu}, \underline{\mathbf{K}})$$

$\underline{K}_{ij} = \underline{\kappa}(\mathbf{x}_i, \mathbf{x}_j)$ and $\underline{\mu} = (m(\mathbf{x}_1), \dots, m(\mathbf{x}_N))^T$

Training

$\mathcal{D} = \{(\mathbf{x}_i, f_i), i = 1 : N\}$, where $f_i = f(\mathbf{x}_i)$ is the noise-free

$$\rightarrow \begin{pmatrix} \mathbf{f} \\ \mathbf{f}_* \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} \underline{\mu} \\ \underline{\mu}_* \end{pmatrix}, \begin{pmatrix} \underline{\mathbf{K}} & \underline{\mathbf{K}}_* \\ \underline{\mathbf{K}}_*^T & \underline{\mathbf{K}}_{**} \end{pmatrix} \right)$$



$\mathbf{K} = \kappa(\mathbf{X}, \mathbf{X})$ is $N \times N$, $\mathbf{K}_* = \kappa(\mathbf{X}, \mathbf{X}_*)$ is $N \times N_*$, and $\mathbf{K}_{**} = \kappa(\mathbf{X}_*, \mathbf{X}_*)$ is $N_* \times N_*$.

Gaussian Process (Murphy ch15)

$$\begin{pmatrix} \mathbf{f} \\ \mathbf{f}_* \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} \boldsymbol{\mu} \\ \boldsymbol{\mu}_* \end{pmatrix}, \begin{pmatrix} \mathbf{K} & \mathbf{K}_* \\ \mathbf{K}_*^T & \mathbf{K}_{**} \end{pmatrix} \right)$$

The conditional is Gaussian:

$$\left. \begin{aligned} \underline{p(\mathbf{f}_* | \mathbf{X}_*, \mathbf{X}, \mathbf{f})} &= \mathcal{N}(\mathbf{f}_* | \boldsymbol{\mu}_*, \boldsymbol{\Sigma}_*) \\ \boldsymbol{\mu}_* &= \boldsymbol{\mu}(\mathbf{X}_*) + \mathbf{K}_*^T \mathbf{K}^{-1} (\mathbf{f} - \boldsymbol{\mu}(\mathbf{X})) \\ \boldsymbol{\Sigma}_* &= \mathbf{K}_{**} - \mathbf{K}_*^T \mathbf{K}^{-1} \mathbf{K}_* \end{aligned} \right\}$$

Common kernel is the squared exponential, RBF, Gaussian kernel

$$\underline{\kappa(x, x')} = \sigma_f^2 \exp\left(-\frac{1}{2\ell^2}(x - x')^2\right)$$

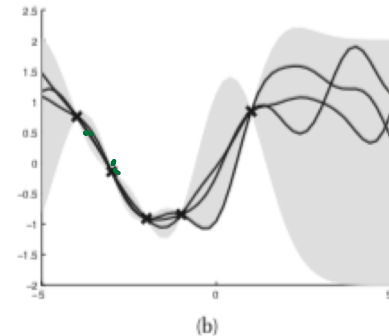
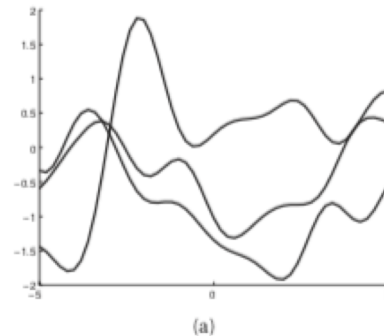


Figure 15.2 Left: some functions sampled from a GP prior with SE kernel. Right: some samples from a GP posterior, after conditioning on 5 noise-free observations. The shaded area represents $\mathbb{E}[f(\mathbf{x})] \pm 2\text{std}(f(\mathbf{x}))$. Based on Figure 2.2 of (Rasmussen and Williams 2006). Figure generated by `gprDemoNoiseFree`.