Announcements

X

Cody homework emailed.

Due Monday and Wednesday before class

Piazza to come Email me, if just attending as that way you can look at homework.

Podcast might work eventually. Introduction to ML, lecture 20 ;-)

Grade last year (A+ 19, A 20, A- 13, B+ 7, S 1, W 1)

Today:

- Gaussian 1.2
- Decision theory 1.5
- Information theory 1.6
- Homework
- Gaussian 2.3

Monday

Gaussian 2.3, Non parametric 1.5, Linear models for regression 3



Likelihood $L(w) = \frac{N}{T} M(z_{n} | w|_{x_{n}}^{2} s^{2}) = \frac{1}{T} \frac{1}{\sqrt{2\pi s^{2}}} e^{z}$ dr = 1 Same differentiation ひっこうなってないない。 Model $t = W^T \times + n$ Maximum Likelihood Bishop 1.2.5 C 2 2(2, -w/2) C Uredient Minizedlon n~ M(0, 5) = N(0, 13-1) t~ M(0, 5) = N(0, 13-1)

Maximum Likelihood

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

(1.61)

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

n = 1

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

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

Giving estimates of W and beta, we can predict

 $eta_{ ext{ML}}$

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

Predictive Distribution





MAP: A Step towards Bayes 1.2.5

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

$$\frac{p(\mathbf{w}|\mathbf{x}, \mathbf{t}, \alpha, \beta) \propto p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta)p(\mathbf{w}|\alpha) = \mathcal{N}\left\{-\frac{\alpha}{2}\mathbf{w}^{\mathrm{T}}\mathbf{w}\right\}$$

$$\frac{p(\mathbf{w}|\mathbf{x}, \mathbf{t}, \alpha, \beta) \propto p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta)p(\mathbf{w}|\alpha) = \mathcal{N}\left\{-\frac{\alpha}{2}\mathbf{w}^{\mathrm{T}}\mathbf{w}\right\}$$

$$\frac{p(\mathbf{w}|\mathbf{x}, \mathbf{t}, \alpha, \beta) \propto p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta)p(\mathbf{w}|\alpha) = \mathcal{N}\left\{-\frac{\alpha}{2}\mathbf{w}^{\mathrm{T}}\mathbf{w}\right\}$$

$$\beta \widetilde{E}(\mathbf{w}) = \frac{\beta}{2} \sum_{n=1}^{N} \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{\alpha}{2} \mathbf{w}^{\mathrm{T}} \mathbf{w}$$

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

Regularized sum of squares

Entropy 1.6

$$H[x] = -\sum_{x} \hat{p}(x) \log_2 p(x)$$

- Important quantity in coding theory statistical physics machine learning





Differential Entropy

Put bins of width ¢ along the real line

$$\lim_{\Delta \to 0} \left\{ -\sum_{i} p(x_i) \Delta \ln p(x_i) \right\} = -\int p(x) \ln p(x) \, \mathrm{d}x$$

For fixed σ^2 differential entropy maximized when

in which case

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

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

The Kullback-Leibler Divergence

P true distribution, q is approximating distribution

$$\operatorname{KL}(p||q) = -\int p(\mathbf{x}) \ln q(\mathbf{x}) \, \mathrm{d}\mathbf{x} - \left(-\int p(\mathbf{x}) \ln p(\mathbf{x}) \, \mathrm{d}\mathbf{x}\right)$$
$$= -\int p(\mathbf{x}) \ln \left\{\frac{q(\mathbf{x})}{p(\mathbf{x})}\right\} \, \mathrm{d}\mathbf{x}$$

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

$$\operatorname{KL}(p\|q) \ge 0$$
 $\operatorname{KL}(p\|q) \not\equiv \operatorname{KL}(q\|p)$

Decision Theory

Inference step

Determine either $p(t|\mathbf{x})$ or $p(\mathbf{x}, t)$.

Decision step

For given x, determine optimal t.





Mixtures of Gaussians (Bishop 2.3.9)

Old Faithful geyser:

than 2 $1/_2$ minutes, or 91 minutes after an eruption lasting more than 2 $1/_2$ minutes. ±10 minutes, Old Faithful will erupt either 65 minutes after an eruption lasting less or 91 minutes, and is dependent on the length of the prior eruption. Within a margin of error of The time between eruptions has a bimodal distribution, with the mean interval being either 65











Mixtures of Gaussians (Bishop 2.3.9)

Determining parameters π , μ , and Σ using maximum log likelihood

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

Solution: use standard, iterative, numeric optimization methods or the expectation maximization algorithm (Chapter 9). Log of a sum; no closed form maximum.

Homework

Recall Curve Fitting





The Gaussian Distribution

Central Limit Theorem

Gaussian as N grows. •The distribution of the sum of N i.i.d. random variables becomes increasingly

•Example: N uniform [0,1] random variables.



Geometry of the Multivariate Gaussian



Moments of the Multivariate Gaussian (2)

$$\mathbb{E}[\mathbf{x}\mathbf{x}^{\mathrm{T}}] = \boldsymbol{\mu}\boldsymbol{\mu}^{\mathrm{T}} + \boldsymbol{\Sigma}$$
$$\operatorname{cov}[\mathbf{x}] = \mathbb{E}\left[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^{\mathrm{T}}\right] = \boldsymbol{\Sigma} = \mathbf{\Sigma}$$

A Gaussian requires D*(D-1)/2 +D parameters. Often we use D +D or Just D+1 parameters.

[]





Maximum Likelihood for the Gaussian (1)
Given i.i.d. data
$$\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)^T$$
, the log likelihood function is given by
 $\ln p(\mathbf{X}|\mu, \Sigma) = -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln|\Sigma| - \frac{1}{2} \sum_{n=1}^{N} (\mathbf{x}_n - \mu)^T \Sigma^{-1} (\mathbf{x}_n - \mu)$
 $\mathcal{Y} = -\frac{1}{2} \ln(2\pi) - \frac{N}{2} \ln|\Sigma| - \frac{1}{2} \sum_{n=1}^{N} (\mathbf{x}_n - \mu)^T \Sigma^{-1} (\mathbf{x}_n - \mu)$
 $\mathcal{Y} = -\frac{1}{2} \ln|\Sigma| + \frac{1}{2} \ln|\Sigma| + \frac{1}{2} \ln \frac{1}{2} \sum_{n=1}^{N} (\mathbf{x}_n - \mu)^T \Sigma^{-1} (\mathbf{x}_n - \mu)$
 $\mathcal{Y} = -\frac{1}{2} \sum_{n=1}^{N} (\mathbf{x}_n - \mu) = 0$
 $\mathcal{Y} = \frac{1}{2} \ln|\Sigma| + \frac{1}{2} \ln|\Sigma| + \frac{1}{2} \ln \frac{1}{2} \sum_{n=1}^{N} (\mathbf{x}_n - \mu)^T \Sigma^{-1} (\mathbf{x}_n - \mu)$
 $\mathcal{Y} = -\frac{1}{2} \sum_{n=1}^{N} (\mathbf{x}_n - \mu) = 0$
 $\mathcal{Y} = \frac{1}{2} \sum_{n=1}^{N} (\mathbf{x}_n - \mu)^T$
 $\mathcal{Y} = -\frac{1}{2} \sum_{n=1}^{N} (\mathbf{x}_n - \mu)^T$

$$\frac{\partial}{\partial \mathbf{A}} \ln |\mathbf{A}| = (\mathbf{A}^{-1})^{\mathrm{T}}$$
$$\frac{\partial}{\partial \mathbf{A}} \operatorname{Tr} (\mathbf{A} \mathbf{B}) = \mathbf{\dot{B}}^{\mathrm{T}} \cdot \mathbf{\dot{A}}$$
$$\frac{\partial}{\partial x} (\mathbf{A}^{-1}) = -\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x} \mathbf{A}^{-1}$$

(C.28) (C.24) (C.21)

Maximum Likelihood for the Gaussian

Set the derivative of the log likelihood function to zero,

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

• Similarly
$$\boldsymbol{\mu}_{\mathrm{ML}} = rac{1}{N}\sum_{n=1}^{N}\mathbf{x}_{n}.$$

$$oldsymbol{\Sigma}_{ ext{ML}} = rac{1}{N} \sum_{n=1}^{N} (\mathbf{x}_n - oldsymbol{\mu}_{ ext{ML}}) (\mathbf{x}_n - oldsymbol{\mu}_{ ext{ML}})^{ ext{T}}.$$