

10-701
Machine Learning
Recitation #3

**MLE, MAP, and Vector/Matrix
Differentiation**

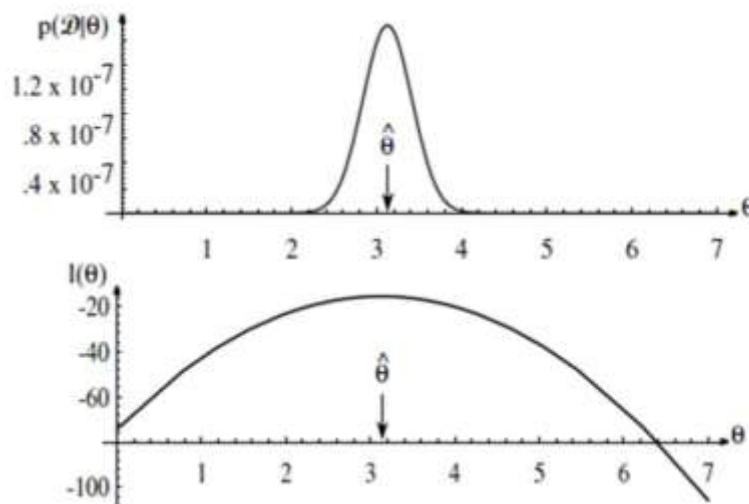
Calvin Murdock
cmurdock@andrew.cmu.edu

(Some slides borrowed from Andrew Moore, Steven Nydick, etc.)

Maximum Likelihood Estimation (MLE)

$$\mathcal{L}(\boldsymbol{\theta}) = \ln p(\mathcal{D}|\boldsymbol{\theta}) = \ln \prod_{i=1}^N p(\mathbf{x}^{(i)}|\boldsymbol{\theta}) = \sum_{i=1}^N \ln p(\mathbf{x}^{(i)}|\boldsymbol{\theta})$$

$$\hat{\boldsymbol{\theta}}_{ML} = \operatorname{argmax}_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) = \operatorname{argmax}_{\boldsymbol{\theta}} \sum_{i=1}^N \ln p(\mathbf{x}^{(i)}|\boldsymbol{\theta})$$



MLE for univariate Gaussian

- Suppose you have $x_1, x_2, \dots, x_R \sim N(\mu, \sigma^2)$
- But you don't know μ

(you do know σ^2)

MLE: For which μ is x_1, x_2, \dots, x_R most likely?

MAP: Which μ maximizes $p(\mu | x_1, x_2, \dots, x_R, \sigma^2)$?

MLE for univariate Gaussian

- Suppose you have $x_1, x_2, \dots, x_R \sim N(\mu, \sigma^2)$
- But you don't know μ (you do know σ^2)
- **MLE:** For which μ is x_1, x_2, \dots, x_R most likely?

$$\mu^{mle} = \arg \max_{\mu} p(x_1, x_2, \dots, x_R \mid \mu, \sigma^2)$$

MLE for univariate Gaussian

$$\mu^{mle} = \arg \max_{\mu} p(x_1, x_2, \dots, x_R \mid \mu, \sigma^2)$$

MLE for univariate Gaussian

$$\mu^{mle} = \arg \max_{\mu} p(x_1, x_2, \dots, x_R \mid \mu, \sigma^2)$$

$$= \arg \max_{\mu} \prod_{i=1}^R p(x_i \mid \mu, \sigma^2) \quad \text{(by i.i.d)}$$

$$= \arg \max_{\mu} \sum_{i=1}^R \log p(x_i \mid \mu, \sigma^2) \quad \text{(monotonicity of log)}$$

$$= \arg \max_{\mu} \frac{1}{\sqrt{2\pi} \sigma} \sum_{i=1}^R -\frac{(x_i - \mu)^2}{2\sigma^2} \quad \text{(plug in formula for Gaussian)}$$

$$= \arg \min_{\mu} \sum_{i=1}^R (x_i - \mu)^2 \quad \text{(after simplification)}$$

Intermission: A General Scalar MLE strategy

Task: Find MLE θ assuming known form for $p(\text{Data} \mid \theta, \text{stuff})$

1. Write $LL = \log P(\text{Data} \mid \theta, \text{stuff})$
2. Work out $\partial LL / \partial \theta$ using high-school calculus
3. Set $\partial LL / \partial \theta = 0$ for a maximum, creating an equation in terms of θ
4. Solve it*
5. Check that you've found a maximum rather than a minimum or saddle-point, and be careful if θ is constrained

*This is a perfect example of something that works perfectly in all textbook examples and usually involves surprising pain if you need it for something new.

The MLE

$$\mu^{mle} = \arg \max_{\mu} p(x_1, x_2, \dots, x_R \mid \mu, \sigma^2)$$

$$= \arg \min_{\mu} \sum_{i=1}^R (x_i - \mu)^2$$

$$= \mu \text{ s.t. } 0 = \frac{\partial \text{LL}}{\partial \mu} =$$

= (what?)

The MLE

$$\mu^{mle} = \arg \max_{\mu} p(x_1, x_2, \dots, x_R \mid \mu, \sigma^2)$$

$$= \arg \min_{\mu} \sum_{i=1}^R (x_i - \mu)^2$$

$$= \mu \text{ s.t. } 0 = \frac{\partial \text{LL}}{\partial \mu} = \frac{\partial}{\partial \mu} \sum_{i=1}^R (x_i - \mu)^2 \\ - \sum_{i=1}^R 2(x_i - \mu)$$

$$\text{Thus } \mu = \frac{1}{R} \sum_{i=1}^R x_i$$

The MLE

$$\mu^{mle} = \frac{1}{R} \sum_{i=1}^R x_i$$

- The best estimate of the mean of a distribution is the mean of the sample!

A General MLE strategy

Suppose $\theta = (\theta_1, \theta_2, \dots, \theta_n)^\top$ is a vector of parameters.

Task: Find MLE θ assuming known form for $p(\text{Data} | \theta, \text{stuff})$

1. Write $LL = \log P(\text{Data} | \theta, \text{stuff})$
2. Work out the gradient $\partial LL / \partial \theta$ using high-school calculus

$$\nabla_{\theta} LL = \frac{\partial LL}{\partial \theta} = \begin{pmatrix} \frac{\partial LL}{\partial \theta_1} \\ \frac{\partial LL}{\partial \theta_2} \\ \vdots \\ \frac{\partial LL}{\partial \theta_n} \end{pmatrix}$$

A General MLE strategy

Suppose $\theta = (\theta_1, \theta_2, \dots, \theta_n)^\top$ is a vector of parameters.

Task: Find MLE θ assuming known form for $p(\text{Data} | \theta, \text{stuff})$

1. Write $LL = \log P(\text{Data} | \theta, \text{stuff})$
2. Work out the gradient $\partial LL / \partial \theta$ using high-school calculus
3. Solve the set of simultaneous equations

$$\begin{aligned}\frac{\partial LL}{\partial \theta_1} &= 0 \\ \frac{\partial LL}{\partial \theta_2} &= 0 \\ &\vdots \\ \frac{\partial LL}{\partial \theta_n} &= 0\end{aligned}$$

A General MLE strategy

Suppose $\theta = (\theta_1, \theta_2, \dots, \theta_n)^\top$ is a vector of parameters.

Task: Find MLE θ assuming known form for $p(\text{Data} | \theta, \text{stuff})$

1. Write $LL = \log P(\text{Data} | \theta, \text{stuff})$
2. Work out the gradient $\partial LL / \partial \theta$ using high-school calculus
3. Solve the set of simultaneous equations

$$\frac{\partial LL}{\partial \theta_1} = 0$$

$$\frac{\partial LL}{\partial \theta_2} = 0$$

$$\vdots$$

$$\frac{\partial LL}{\partial \theta_n} = 0$$

4. Check that you're at a maximum

A General MLE strategy

Suppose $\theta = (\theta_1, \theta_2, \dots, \theta_n)^\top$ is a vector of parameters.

Task: Find MLE θ assuming known form for $p(\text{Data} | \theta, \text{stuff})$

1. Write $LL = \log P(\text{Data} | \theta, \text{stuff})$
2. Work out the gradient $\partial LL / \partial \theta$ using high-school calculus
3. Solve the set of simultaneous equations

If you can't solve them,
what should you do?

$$\begin{aligned}\frac{\partial LL}{\partial \theta_1} &= 0 \\ \frac{\partial LL}{\partial \theta_2} &= 0 \\ &\vdots \\ \frac{\partial LL}{\partial \theta_n} &= 0\end{aligned}$$

4. Check that you're at a maximum

MLE for univariate Gaussian

- Suppose you have $x_1, x_2, \dots, x_R \sim (\text{i.i.d}) \mathcal{N}(\mu, \sigma^2)$
- But you don't know μ or σ^2
- MLE: For which $\theta = (\mu, \sigma^2)$ is x_1, x_2, \dots, x_R most likely?

$$\log p(x_1, x_2, \dots, x_R \mid \mu, \sigma^2) = -R(\log \pi + \frac{1}{2} \log \sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^R (x_i - \mu)^2$$

$$\frac{\partial LL}{\partial \mu} = \frac{1}{\sigma^2} \sum_{i=1}^R (x_i - \mu)$$

$$\frac{\partial LL}{\partial \sigma^2} = -\frac{R}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^R (x_i - \mu)^2$$

MLE for univariate Gaussian

- Suppose you have $x_1, x_2, \dots, x_R \sim (\text{i.i.d}) \mathcal{N}(\mu, \sigma^2)$
- But you don't know μ or σ^2
- MLE: For which $\theta = (\mu, \sigma^2)$ is x_1, x_2, \dots, x_R most likely?

$$\log p(x_1, x_2, \dots, x_R \mid \mu, \sigma^2) = -R(\log \pi + \frac{1}{2} \log \sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^R (x_i - \mu)^2$$

$$0 = \frac{1}{\sigma^2} \sum_{i=1}^R (x_i - \mu)$$

$$0 = -\frac{R}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^R (x_i - \mu)^2$$

MLE for univariate Gaussian

- Suppose you have $x_1, x_2, \dots, x_R \sim (\text{i.i.d}) \mathcal{N}(\mu, \sigma^2)$
- But you don't know μ or σ^2
- MLE: For which $\theta = (\mu, \sigma^2)$ is x_1, x_2, \dots, x_R most likely?

$$\log p(x_1, x_2, \dots, x_R \mid \mu, \sigma^2) = -R(\log \pi + \frac{1}{2} \log \sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^R (x_i - \mu)^2$$

$$0 = \frac{1}{\sigma^2} \sum_{i=1}^R (x_i - \mu) \Rightarrow \mu = \frac{1}{R} \sum_{i=1}^R x_i$$

$$0 = -\frac{R}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^R (x_i - \mu)^2 \Rightarrow \text{what?}$$

MLE for univariate Gaussian

- Suppose you have $x_1, x_2, \dots, x_R \sim (\text{i.i.d}) \mathcal{N}(\mu, \sigma^2)$
- But you don't know μ or σ^2
- MLE: For which $\theta = (\mu, \sigma^2)$ is x_1, x_2, \dots, x_R most likely?

$$\mu^{mle} = \frac{1}{R} \sum_{i=1}^R x_i$$

$$\sigma_{mle}^2 = \frac{1}{R} \sum_{i=1}^R (x_i - \mu^{mle})^2$$

Unbiased Estimators

- An estimator of a parameter is **unbiased** if the expected value of the estimate is the **same** as the true value of the parameters.
- If $x_1, x_2, \dots, x_R \sim (\text{i.i.d}) N(\mu, \sigma^2)$ then

$$E[\mu^{mle}] = E\left[\frac{1}{R} \sum_{i=1}^R x_i\right] = \mu$$

μ^{mle} is unbiased

Biased Estimators

- An estimator of a parameter is **biased** if the expected value of the estimate is **different from** the true value of the parameters.
- If $x_1, x_2, \dots, x_R \sim (\text{i.i.d}) N(\mu, \sigma^2)$ then

$$E[\sigma_{mle}^2] = E\left[\frac{1}{R} \sum_{i=1}^R (x_i - \mu^{mle})^2\right] = E\left[\frac{1}{R} \left(\sum_{i=1}^R x_i - \frac{1}{R} \sum_{j=1}^R x_j\right)^2\right] \neq \sigma^2$$

σ_{mle}^2 is biased

MLE Variance Bias

- If $x_1, x_2, \dots, x_R \sim (\text{i.i.d}) N(\mu, \sigma^2)$ then

$$E[\sigma_{mle}^2] = E\left[\frac{1}{R}\left(\sum_{i=1}^R x_i - \frac{1}{R}\sum_{j=1}^R x_j\right)^2\right] = \left(1 - \frac{1}{R}\right)\sigma^2 \neq \sigma^2$$

Intuition check: consider the case of $R=1$

Why should our guts expect that σ_{mle}^2 would be an underestimate of true σ^2 ?

How could you prove that?

Unbiased estimate of Variance

- If $x_1, x_2, \dots, x_R \sim (\text{i.i.d}) N(\mu, \sigma^2)$ then

$$E[\sigma_{mle}^2] = E\left[\frac{1}{R}\left(\sum_{i=1}^R x_i - \frac{1}{R}\sum_{j=1}^R x_j\right)^2\right] = \left(1 - \frac{1}{R}\right)\sigma^2 \neq \sigma^2$$

So define $\sigma_{\text{unbiased}}^2 = \frac{\sigma_{mle}^2}{\left(1 - \frac{1}{R}\right)}$ So $E[\sigma_{\text{unbiased}}^2] = \sigma^2$

Unbiased estimate of Variance

- If $x_1, x_2, \dots, x_R \sim (\text{i.i.d}) N(\mu, \sigma^2)$ then

$$E[\sigma_{mle}^2] = E\left[\frac{1}{R}\left(\sum_{i=1}^R x_i - \frac{1}{R}\sum_{j=1}^R x_j\right)^2\right] = \left(1 - \frac{1}{R}\right)\sigma^2 \neq \sigma^2$$

So define $\sigma_{\text{unbiased}}^2 = \frac{\sigma_{mle}^2}{\left(1 - \frac{1}{R}\right)}$ So $E[\sigma_{\text{unbiased}}^2] = \sigma^2$

$$\sigma_{\text{unbiased}}^2 = \frac{1}{R-1} \sum_{i=1}^R (x_i - \mu^{mle})^2$$

Unbiaseditude discussion

- *Which is best?*

$$\sigma_{mle}^2 = \frac{1}{R} \sum_{i=1}^R (x_i - \mu^{mle})^2$$

$$\sigma_{\text{unbiased}}^2 = \frac{1}{R-1} \sum_{i=1}^R (x_i - \mu^{mle})^2$$

Answer:

- It depends on the task
- And doesn't make much difference once $R \rightarrow \text{large}$

Don't get too excited about being unbiased

- *Assume* $x_1, x_2, \dots, x_R \sim (\text{i.i.d}) \mathcal{N}(\mu, \sigma^2)$
- Suppose we had these estimators for the mean

$$\mu^{\text{suboptimal}} = \frac{1}{R + 7\sqrt{R}} \sum_{i=1}^R x_i$$

$$\mu^{\text{crap}} = x_1$$

Are either of these unbiased?

Will either of them asymptote to the correct value as R gets large?

Which is more useful?

Maximum Conditional Likelihood Estimation (MCLE)

- Same as MLE except with conditional likelihood.
 - E.g. for regression: Given $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$, learn $\boldsymbol{\theta}$

Diagram illustrating the regression model and its equivalent probabilistic representation:

Output y is determined by the Given input \mathbf{x} through a Deterministic function with parameters $\boldsymbol{\theta}$, plus Gaussian noise ϵ .

$$y = f_{\boldsymbol{\theta}}(\mathbf{x}) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

Equivalent to:

$$y \sim \mathcal{N}(f_{\boldsymbol{\theta}}(\mathbf{x}), \sigma^2)$$

where the mean μ is $f_{\boldsymbol{\theta}}(\mathbf{x})$.

Maximum Conditional Likelihood Estimation (MCLE)

- Same as MLE except with conditional likelihood.
 - E.g. for regression: Given $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$, learn θ

Output Given input Gaussian noise

$y = f_{\theta}(\mathbf{x}) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$

Deterministic function with parameters θ Equivalent to:

$y \sim \mathcal{N}(f_{\theta}(\mathbf{x}), \sigma^2)$

Standard MLE: $\mu^{mle} = \arg \max_{\mu} p(y_1, \dots, y_n | \mu, \sigma^2)$

Maximum Conditional Likelihood Estimation (MCLE)

- Same as MLE except with conditional likelihood.
 - E.g. for regression: Given $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$, learn $\boldsymbol{\theta}$

Diagram illustrating the regression model and its equivalent probabilistic representation:

Output y is determined by the Given input \mathbf{x} through a Deterministic function with parameters $\boldsymbol{\theta}$, plus Gaussian noise ϵ .

$$y = f_{\boldsymbol{\theta}}(\mathbf{x}) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

Equivalent to:

$$y \sim \mathcal{N}(f_{\boldsymbol{\theta}}(\mathbf{x}), \sigma^2)$$

where the mean $\mu = f_{\boldsymbol{\theta}}(\mathbf{x})$.

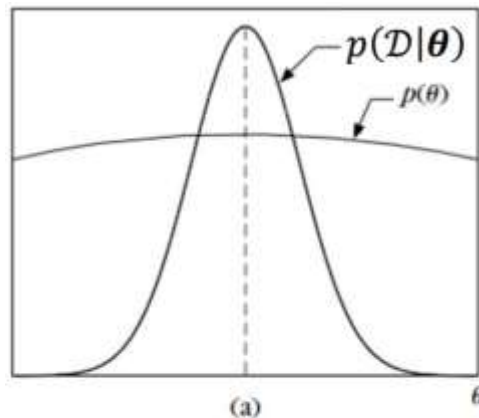
Conditional MLE: $\boldsymbol{\theta}^{mle} = \arg \max_{\boldsymbol{\theta}} p(y_1, \dots, y_n | \boldsymbol{\theta}, \sigma^2, \mathbf{x}_1, \dots, \mathbf{x}_n)$

Maximum a Posteriori (MAP) Estimation

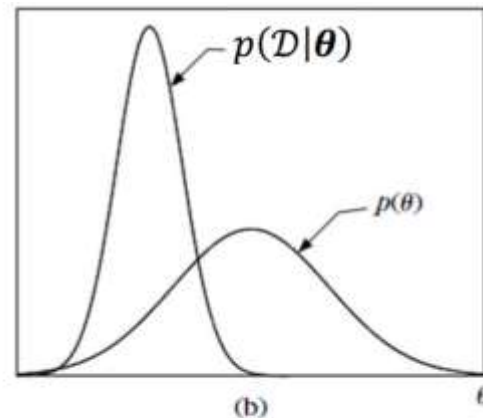
$$\hat{\boldsymbol{\theta}}_{MAP} = \operatorname{argmax}_{\boldsymbol{\theta}} p(\boldsymbol{\theta}|\mathcal{D})$$

- ▶ Since $p(\boldsymbol{\theta}|\mathcal{D}) \propto p(\mathcal{D}|\boldsymbol{\theta})p(\boldsymbol{\theta})$

$$\hat{\boldsymbol{\theta}}_{MAP} = \operatorname{argmax}_{\boldsymbol{\theta}} p(\mathcal{D}|\boldsymbol{\theta})p(\boldsymbol{\theta})$$



$$\hat{\boldsymbol{\theta}}_{MAP} \cong \hat{\boldsymbol{\theta}}_{ML}$$



$$\hat{\boldsymbol{\theta}}_{MAP} > \hat{\boldsymbol{\theta}}_{ML}$$

Being Bayesian: MAP estimates for Gaussians

- Suppose you have $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_R \sim (\text{i.i.d}) \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$
- But you don't know $\boldsymbol{\mu}$ or $\boldsymbol{\Sigma}$
- MAP: Which $(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ maximizes $p(\boldsymbol{\mu}, \boldsymbol{\Sigma} \mid \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_R)$?

Step 1: Put a prior on $(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

Being Bayesian: MAP estimates for Gaussians

- Suppose you have $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_R \sim (\text{i.i.d}) \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$
- But you don't know $\boldsymbol{\mu}$ or $\boldsymbol{\Sigma}$
- MAP: Which $(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ maximizes $p(\boldsymbol{\mu}, \boldsymbol{\Sigma} \mid \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_R)$?

Step 1: Put a prior on $(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

Step 1a: Put a prior on $\boldsymbol{\Sigma}$

$$(\nu_0 - m - 1) \boldsymbol{\Sigma} \sim \text{IW}(\nu_0, (\nu_0 - m - 1) \boldsymbol{\Sigma}_0)$$

This thing is called the Inverse-Wishart distribution.

A PDF over SPD matrices!

Being Bayesian: MAP estimates for Gaussians

- Suppose you have $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_R \sim (\text{i.i.d}) \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$
- But you don't know $\boldsymbol{\mu}$ or $\boldsymbol{\Sigma}$
- MAP: Which $(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ maximizes $p(\boldsymbol{\mu}, \boldsymbol{\Sigma} \mid \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_R)$?

Step 1: Put a prior on $(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

Step 1a: Put a prior on $\boldsymbol{\Sigma}$

$$(\nu_0 - m - 1) \boldsymbol{\Sigma} \sim \text{IW}(\nu_0, (\nu_0 - m - 1) \boldsymbol{\Sigma}_0)$$

This thing is called the Inverse-Wishart distribution.

A PDF over SPD matrices!

Being Bayesian: MAP estimates for Gaussians

- Suppose you have $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_R \sim (\text{i.i.d}) \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$
- But you don't know $\boldsymbol{\mu}$ or $\boldsymbol{\Sigma}$
- MAP: Which $(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ maximizes $p(\boldsymbol{\mu}, \boldsymbol{\Sigma} \mid \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_R)$?

Step 1: Put a prior on $(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

Step 1a: Put a prior on $\boldsymbol{\Sigma}$

$$(\nu_0 - m - 1)\boldsymbol{\Sigma} \sim \text{IW}(\nu_0, (\nu_0 - m - 1)\boldsymbol{\Sigma}_0)$$

Step 1b: Put a prior on $\boldsymbol{\mu} \mid \boldsymbol{\Sigma}$

$$\boldsymbol{\mu} \mid \boldsymbol{\Sigma} \sim \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma} / \kappa_0)$$

Together, " $\boldsymbol{\Sigma}$ " and " $\boldsymbol{\mu} \mid \boldsymbol{\Sigma}$ " define a joint distribution on $(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

Being Bayesian: MAP estimates for Gaussians

- Suppose you have $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_R \sim (\text{i.i.d}) \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$
- But you don't know $\boldsymbol{\mu}$ or $\boldsymbol{\Sigma}$
- MAP: Which $(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ maximizes $p(\boldsymbol{\mu}, \boldsymbol{\Sigma} \mid \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_R)$?

Step 1: Put a prior on $(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

Why do we use this form of prior?

Step 1a: Put a prior on $\boldsymbol{\Sigma}$

$$(\nu_0 - m - 1)\boldsymbol{\Sigma} \sim \text{IW}(\nu_0, (\nu_0 - m - 1)\boldsymbol{\Sigma}_0)$$

Step 1b: Put a prior on $\boldsymbol{\mu} \mid \boldsymbol{\Sigma}$

$$\boldsymbol{\mu} \mid \boldsymbol{\Sigma} \sim \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma} / \kappa_0)$$

Being Bayesian: MAP estimates for Gaussians

- Suppose you have $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_R \sim (\text{i.i.d}) \mathcal{N}(\boldsymbol{\mu}, \Sigma)$
- But you don't know $\boldsymbol{\mu}$ or Σ
- MAP: Which $(\boldsymbol{\mu}, \Sigma)$ maximizes $p(\boldsymbol{\mu}, \Sigma \mid \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_R)$?

Step 1: Put a prior on $(\boldsymbol{\mu}, \Sigma)$

Step 1a: Put a prior on Σ

$$(\nu_0 - m - 1)\Sigma \sim \text{IW}(\nu_0, (\nu_0 - m - 1)\Sigma_0)$$

Step 1b: Put a prior on $\boldsymbol{\mu} \mid \Sigma$

$$\boldsymbol{\mu} \mid \Sigma \sim \mathcal{N}(\boldsymbol{\mu}_0, \Sigma / \kappa_0)$$

Why do we use this form of prior?

Actually, we don't have to

But it is computationally and algebraically convenient...

...it's a *conjugate prior*.

Being Bayesian: MAP estimates for Gaussians

- Suppose you have $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_R \sim (\text{i.i.d}) \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$
- MAP: Which $(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ maximizes $p(\boldsymbol{\mu}, \boldsymbol{\Sigma} \mid \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_R)$?

Step 1: Prior: $(\nu_0 - m - 1) \boldsymbol{\Sigma} \sim \text{IW}(\nu_0, (\nu_0 - m - 1) \boldsymbol{\Sigma}_0)$, $\boldsymbol{\mu} \mid \boldsymbol{\Sigma} \sim \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma} / \kappa_0)$

Step 2:

$$\bar{\mathbf{x}} = \frac{1}{R} \sum_{k=1}^R \mathbf{x}_k$$

$$\boldsymbol{\mu}_R = \frac{\kappa_0 \boldsymbol{\mu}_0 + R \bar{\mathbf{x}}}{\kappa_0 + R}$$

$$\nu_R = \nu_0 + R$$

$$\kappa_R = \kappa_0 + R$$

$$(\nu_R + m - 1) \boldsymbol{\Sigma}_R = (\nu_0 + m - 1) \boldsymbol{\Sigma}_0 + \sum_{k=1}^R (\mathbf{x}_k - \bar{\mathbf{x}})(\mathbf{x}_k - \bar{\mathbf{x}})^T + \frac{(\bar{\mathbf{x}} - \boldsymbol{\mu}_0)(\bar{\mathbf{x}} - \boldsymbol{\mu}_0)^T}{1/\kappa_0 + 1/R}$$

Step 3: Posterior: $(\nu_R + m - 1) \boldsymbol{\Sigma} \sim \text{IW}(\nu_R, (\nu_R + m - 1) \boldsymbol{\Sigma}_R)$,

$$\boldsymbol{\mu} \mid \boldsymbol{\Sigma} \sim \mathcal{N}(\boldsymbol{\mu}_R, \boldsymbol{\Sigma} / \kappa_R)$$

Result: $\boldsymbol{\mu}^{\text{map}} = \boldsymbol{\mu}_R$, $E[\boldsymbol{\Sigma} \mid \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_R] = \boldsymbol{\Sigma}_R$

Being Bayesian: M

- Suppose you have \mathbf{x}
- MAP: Which $(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ m

- Look carefully at what these formulae are doing. It's all very sensible.
- Conjugate priors mean prior form and posterior form are same and characterized by "sufficient statistics" of the data.
- The marginal distribution on $\boldsymbol{\mu}$ is a student-t
- One point of view: it's pretty academic if $R > 30$

Step 1: Prior: $(\nu_0 - m - 1) \boldsymbol{\Sigma} \sim$

Step 2:

$$\bar{\mathbf{x}} = \frac{1}{R} \sum_{k=1}^R \mathbf{x}_k$$

$$\boldsymbol{\mu}_R = \frac{\kappa_0 \boldsymbol{\mu}_0 + R \bar{\mathbf{x}}}{\kappa_0 + R}$$

$$\nu_R = \nu_0 + R$$

$$\kappa_R = \kappa_0 + R$$

$$(\nu_R + m - 1) \boldsymbol{\Sigma}_R = (\nu_0 + m - 1) \boldsymbol{\Sigma}_0 + \sum_{k=1}^R (\mathbf{x}_k - \bar{\mathbf{x}})(\mathbf{x}_k - \bar{\mathbf{x}})^T + \frac{(\bar{\mathbf{x}} - \boldsymbol{\mu}_0)(\bar{\mathbf{x}} - \boldsymbol{\mu}_0)^T}{1/\kappa_0 + 1/R}$$

Step 3: Posterior: $(\nu_R + m - 1) \boldsymbol{\Sigma} \sim \text{IW}(\nu_R, (\nu_R + m - 1) \boldsymbol{\Sigma}_R),$

$$\boldsymbol{\mu} \mid \boldsymbol{\Sigma} \sim \text{N}(\boldsymbol{\mu}_R, \boldsymbol{\Sigma} / \kappa_R)$$

Result: $\boldsymbol{\mu}^{\text{map}} = \boldsymbol{\mu}_R, E[\boldsymbol{\Sigma} \mid \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_R] = \boldsymbol{\Sigma}_R$

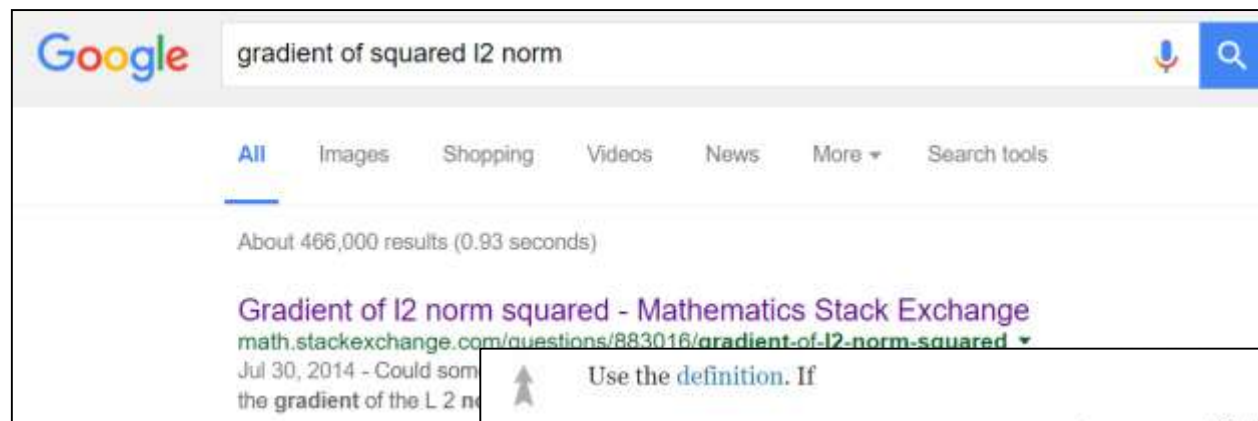
Vector/Matrix Derivatives

- Best reference: **The Matrix Cookbook**
(www.imm.dtu.dk/pubdb/views/edoc_download.php/3274/pdf/imm3274.pdf)

Contents	
1 Basics	6
1.1 Trace	6
1.2 Determinant	6
1.3 The Special Case 2x2	7
2 Derivatives	8
2.1 Derivatives of a Determinant	8
2.2 Derivatives of an Inverse	9
2.3 Derivatives of Eigenvalues	10
2.4 Derivatives of Matrices, Vectors and Scalar Forms	10
2.5 Derivatives of Traces	12
2.6 Derivatives of vector norms	14
2.7 Derivatives of matrix norms	14
2.8 Derivatives of Structured Matrices	14

Vector/Matrix Derivatives

- Second best reference: **Google**



↑ 22
↓
✓

Use the definition. If

$$f(x) = \|x\|_2^2 = \left(\left(\sum_{k=1}^n x_k^2 \right)^{1/2} \right)^2 = \sum_{k=1}^n x_k^2,$$

then

$$\frac{\partial}{\partial x_j} f(x) = \frac{\partial}{\partial x_j} \sum_{k=1}^n x_k^2 = \sum_{k=1}^n \underbrace{\frac{\partial}{\partial x_j} x_k^2}_{\substack{=0, \text{ if } j \neq k, \\ =2x_j, \text{ else}}} = 2x_j.$$


It follows that

$$\nabla f(x) = 2x.$$

share cite improve this answer

edited Jul 30 '14 at 21:00

answered Jul 30 '14 at 20:52

 Surb
27.3k ● 8 ■ 30 ▲ 58

Vector Gradient

A Gradient is the derivative of a scalar with respect to a vector.

scalar function

$$\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} = \left(\left[\frac{\partial f(\mathbf{x})}{\partial x_1} \right] \quad \left[\frac{\partial f(\mathbf{x})}{\partial x_2} \right] \quad \cdots \quad \left[\frac{\partial f(\mathbf{x})}{\partial x_n} \right] \right)^T$$

parameter vector

If we have the function: $f(\mathbf{x}) = 2x_1x_2 + x_2^2 + x_1x_3^2$, then the Gradient is

$$\begin{aligned} \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} &= \left(\left[\frac{\partial f(\mathbf{x})}{\partial x_1} \right] \quad \left[\frac{\partial f(\mathbf{x})}{\partial x_2} \right] \quad \left[\frac{\partial f(\mathbf{x})}{\partial x_3} \right] \right)^T \\ &= [2x_2 + x_3^2 \quad 2x_1 + 2x_2 \quad 2x_1x_3]^T \end{aligned}$$

Matrix Gradient

scalar function

parameter matrix

$$\frac{\partial f(\mathbf{X})}{\partial \mathbf{X}} = \begin{pmatrix} \left[\frac{\partial f(\mathbf{X})}{\partial x_{1,1}} \right] & \cdots & \left[\frac{\partial f(\mathbf{X})}{\partial x_{1,m}} \right] \\ \vdots & \ddots & \vdots \\ \left[\frac{\partial f(\mathbf{X})}{\partial x_{n,1}} \right] & \cdots & \left[\frac{\partial f(\mathbf{X})}{\partial x_{n,m}} \right] \end{pmatrix}$$

Jacobian

A Jacobian is the derivative of a vector with respect to a transposed vector.

vector function
parameter vector

$$\frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}^T} = \begin{pmatrix} \left[\frac{\partial f_1(\mathbf{x})}{\partial x_1} \right] & \cdots & \left[\frac{\partial f_1(\mathbf{x})}{\partial x_n} \right] \\ \vdots & \cdots & \vdots \\ \left[\frac{\partial f_k(\mathbf{x})}{\partial x_1} \right] & \cdots & \left[\frac{\partial f_k(\mathbf{x})}{\partial x_n} \right] \end{pmatrix} = \begin{pmatrix} \left[\frac{\partial f_1(\mathbf{x})}{\partial \mathbf{x}} \right]^T \\ \vdots \\ \left[\frac{\partial f_k(\mathbf{x})}{\partial \mathbf{x}} \right]^T \end{pmatrix}$$

\uparrow
 ∇f_k

If we have the function

$$\mathbf{f}(\mathbf{x}) = [3x_1^2 + x_2 \quad \ln(x_1) \quad \sin(x_2)]^T$$

Then the Jacobian is

$$\frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}^T} = \begin{pmatrix} 6x_1 & 1 \\ \frac{1}{x_1} & 0 \\ 0 & \cos(x_2) \end{pmatrix}$$

Hessian

The Hessian is derivative of a Gradient with respect to a transposed vector.

scalar function
parameter vector

$$\frac{\partial^2 f(\mathbf{x})}{\partial \mathbf{x} \partial \mathbf{x}^T} = \begin{pmatrix} \left[\frac{\partial f(\mathbf{x})}{\partial x_1^2} \right] & \cdots & \left[\frac{\partial f(\mathbf{x})}{\partial x_1 \partial x_n} \right] \\ \vdots & \ddots & \vdots \\ \left[\frac{\partial f(\mathbf{x})}{\partial x_n \partial x_1} \right] & \cdots & \left[\frac{\partial f(\mathbf{x})}{\partial x_n^2} \right] \end{pmatrix}$$

Because our above Gradient is

$$\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} = [2x_2 + x_3^2 \quad 2x_1 + 2x_2 \quad 2x_1x_3]^T$$

The Hessian would be

$$\frac{\partial^2 f(\mathbf{x})}{\partial \mathbf{x} \partial \mathbf{x}^T} = \begin{pmatrix} 0 & 2 & 2x_3 \\ 2 & 2 & 0 \\ 2x_3 & 0 & 2x_1 \end{pmatrix}$$

There are two ways of computing gradients...

$$f(\boldsymbol{\theta}) = \|\mathbf{x} - \boldsymbol{\theta}\|^2 \quad \frac{df}{d\boldsymbol{\theta}} = ?$$

1. Compute each element of gradient using scalar partial derivatives:

$$f(\boldsymbol{\theta}) = \sum_j (x_j - \theta_j)^2 \quad \left(\frac{df}{d\boldsymbol{\theta}}\right)_j = \frac{\partial f}{\partial \theta_j} = \frac{\partial}{\partial \theta_j} (x_j - \theta_j)^2 = 2(\theta_j - x_j)$$

(using chain rule,
power rule)

2. Directly use properties of vector calculus.

$$f(\boldsymbol{\theta}) = (\mathbf{x} - \boldsymbol{\theta})^T (\mathbf{x} - \boldsymbol{\theta}) \quad \frac{df}{d\boldsymbol{\theta}} = 2(\boldsymbol{\theta} - \mathbf{x})$$

(using **vector** chain rule,
vector power rule)

Important Properties

- **Linearity:**

$$\frac{d}{d\mathbf{x}} [a \cdot f(\mathbf{x}) + b \cdot g(\mathbf{x})] = a \cdot \frac{df}{d\mathbf{x}} + b \cdot \frac{dg}{d\mathbf{x}}$$

- **Product Rule:**

$$\frac{d}{d\mathbf{x}} f(\mathbf{x}) \cdot g(\mathbf{x}) = f(\mathbf{x}) \cdot \frac{dg}{d\mathbf{x}} + g(\mathbf{x}) \cdot \frac{df}{d\mathbf{x}}$$

$$\frac{d}{d\mathbf{x}} \mathbf{f}(\mathbf{x})^T \mathbf{g}(\mathbf{x}) = \left(\frac{d\mathbf{g}}{d\mathbf{x}^T} \right)^T \mathbf{f}(\mathbf{x}) + \left(\frac{d\mathbf{f}}{d\mathbf{x}^T} \right)^T \mathbf{g}(\mathbf{x})$$

 Jacobian

- **Chain Rule:**

$$\frac{df}{d\mathbf{x}} = \frac{df}{dg} \cdot \frac{dg}{d\mathbf{x}} = f'(g(\mathbf{x})) \cdot \frac{dg}{d\mathbf{x}}$$

$$f: \mathbb{R} \rightarrow \mathbb{R}, g: \mathbb{R}^p \rightarrow \mathbb{R}$$

$$\frac{df}{d\mathbf{x}} = \left(\frac{d\mathbf{g}}{d\mathbf{x}^T} \right)^T \frac{df}{d\mathbf{g}}$$

$$f: \mathbb{R}^q \rightarrow \mathbb{R}, \mathbf{g}: \mathbb{R}^p \rightarrow \mathbb{R}^q$$

(p × 1) (p × q) (q × 1)

Hint: the sizes should match up.

Examples

$$L = \mathbf{x}^T \mathbf{A} \mathbf{x} = \mathbf{f}(\mathbf{x})^T \mathbf{g}(\mathbf{x})$$

$$\mathbf{f}(\mathbf{x}) = \mathbf{x}, \quad \mathbf{g}(\mathbf{x}) = \mathbf{A} \mathbf{x}$$

$$\frac{dL}{d\mathbf{x}} = \left(\frac{d\mathbf{f}}{d\mathbf{x}^T} \right)^T \mathbf{g}(\mathbf{x}) + \left(\frac{d\mathbf{g}}{d\mathbf{x}^T} \right)^T \mathbf{f}(\mathbf{x})$$

$$\frac{dL}{d\mathbf{x}} = \mathbf{A} \mathbf{x} + \mathbf{A}^T \mathbf{x} = (\mathbf{A} + \mathbf{A}^T) \mathbf{x}$$

Examples

$$L = \mathbf{x}^T \mathbf{A} \mathbf{x} = \mathbf{f}(\mathbf{x})^T \mathbf{g}(\mathbf{x})$$

$$\mathbf{f}(\mathbf{x}) = \mathbf{x}, \quad \mathbf{g}(\mathbf{x}) = \mathbf{A} \mathbf{x}$$

$$\frac{dL}{d\mathbf{x}} = \left(\frac{d\mathbf{f}}{d\mathbf{x}^T} \right)^T \mathbf{g}(\mathbf{x}) + \left(\frac{d\mathbf{g}}{d\mathbf{x}^T} \right)^T \mathbf{f}(\mathbf{x})$$

$$\frac{dL}{d\mathbf{x}} = \mathbf{A} \mathbf{x} + \mathbf{A}^T \mathbf{x} = (\mathbf{A} + \mathbf{A}^T) \mathbf{x}$$

$$L = f(\mathbf{a}^T \mathbf{x}) = f(g(\mathbf{x}))$$

$$g(\mathbf{x}) = \mathbf{a}^T \mathbf{x}$$

$$\frac{dL}{d\mathbf{x}} = \frac{df}{dg} \cdot \frac{dg}{d\mathbf{x}} = f'(g) \cdot \mathbf{a} = f'(\mathbf{a}^T \mathbf{x}) \cdot \mathbf{a}$$

Other Important Properties

- Refer to **The Matrix Cookbook**

$$\frac{\partial \mathbf{x}^T \mathbf{B} \mathbf{x}}{\partial \mathbf{x}} = (\mathbf{B} + \mathbf{B}^T) \mathbf{x} \quad (81)$$

$$\frac{\partial \mathbf{b}^T \mathbf{X}^T \mathbf{D} \mathbf{X} \mathbf{c}}{\partial \mathbf{X}} = \mathbf{D}^T \mathbf{X} \mathbf{b} \mathbf{c}^T + \mathbf{D} \mathbf{X} \mathbf{c} \mathbf{b}^T \quad (82)$$

$$\frac{\partial}{\partial \mathbf{X}} (\mathbf{X} \mathbf{b} + \mathbf{c})^T \mathbf{D} (\mathbf{X} \mathbf{b} + \mathbf{c}) = (\mathbf{D} + \mathbf{D}^T) (\mathbf{X} \mathbf{b} + \mathbf{c}) \mathbf{b}^T \quad (83)$$

Assume \mathbf{W} is symmetric, then

$$\frac{\partial}{\partial \mathbf{s}} (\mathbf{x} - \mathbf{A} \mathbf{s})^T \mathbf{W} (\mathbf{x} - \mathbf{A} \mathbf{s}) = -2 \mathbf{A}^T \mathbf{W} (\mathbf{x} - \mathbf{A} \mathbf{s}) \quad (84)$$

$$\frac{\partial}{\partial \mathbf{x}} (\mathbf{x} - \mathbf{s})^T \mathbf{W} (\mathbf{x} - \mathbf{s}) = 2 \mathbf{W} (\mathbf{x} - \mathbf{s}) \quad (85)$$

- Refer to Wikipedia:
 - <https://en.wikipedia.org/wiki/Gradient>
 - https://en.wikipedia.org/wiki/Matrix_calculus

Next week: convexity, gradient descent, etc...