

Learning with Maximum Likelihood

Andrew W. Moore

Professor

School of Computer Science

Carnegie Mellon University

www.cs.cmu.edu/~awm

awm@cs.cmu.edu

412-268-7599

Note to other teachers and users of these slides. Andrew would be delighted if you found this source material useful in giving your own lectures. Feel free to use these slides verbatim, or to modify them to fit your own needs. PowerPoint originals are available. If you make use of a significant portion of these slides in your own lecture, please include this message, or the following link to the source repository of Andrew's tutorials: <http://www.cs.cmu.edu/~awm/tutorials> . Comments and corrections gratefully received.

Maximum Likelihood learning of Gaussians for Data Mining

- Why we should care
- Learning Univariate Gaussians
- Learning Multivariate Gaussians
- What's a biased estimator?
- Bayesian Learning of Gaussians

Why we should care

- Maximum Likelihood Estimation is a very very very very fundamental part of data analysis.
- “MLE for Gaussians” is training wheels for our future techniques
- Learning Gaussians is more useful than you might guess...

Learning Gaussians from Data

- 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

(you do know σ^2)

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

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

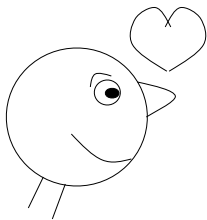
Learning Gaussians from Data

- Suppose you have $x_1, x_2, \dots, x_R \sim (\text{i.i.d}) N(\mu, \sigma^2)$
- 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 $p(\mu|x_1, x_2, \dots, x_R, \sigma^2)$?

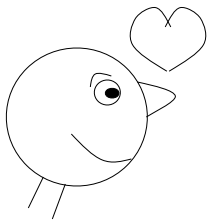
Learning Gaussians from Data

- Suppose you have $x_1, x_2, \dots, x_R \sim (\text{i.i.d}) N(\mu, \sigma^2)$
- 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 $p(\mu|x_1, x_2, \dots, x_R, \sigma^2)$?

Despite this, we'll spend 95% of our time on MLE. Why? Wait and see...

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 (you do know σ^2)
- MLE: For which μ is x_1, x_2, \dots, x_R most likely?

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

Algebra Euphoria

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

=

(by i.i.d)

=

(monotonicity of
log)

=

(plug in formula
for Gaussian)

=

(after
simplification)

Algebra Euphoria

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

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

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

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

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

Intermission: A General Scalar MLE strategy

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

1. Write $LL = \log P(\text{Data} | \theta, \text{stuff})$
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 μ

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

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

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

= (what?)

The MLE μ

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

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

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

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

Lawks-a-lawdy!

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

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

At first sight:

This kind of pedantic, algebra-filled and ultimately unsurprising fact is exactly the reason people throw down their “Statistics” book and pick up their “Agent Based Evolutionary Data Mining Using The Neuro-Fuzz Transform” book.

A General MLE strategy

Suppose $\theta = (\theta_1, \theta_2, \dots, \theta_n)^T$ 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 $\partial LL / \partial \theta$ using high-school calculus

$$\frac{\partial LL}{\partial \theta} = \begin{pmatrix} \frac{\partial LL}{\partial \theta_1} \\ \frac{\partial LL}{\partial \theta_2} \\ \boxed{?} \\ \frac{\partial LL}{\partial \theta_n} \end{pmatrix}$$

A General MLE strategy

Suppose $\theta = (\theta_1, \theta_2, \dots, \theta_n)^T$ 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 $\partial LL / \partial \theta$ using high-school calculus
3. Solve the set of simultaneous equations

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

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

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

A General MLE strategy

Suppose $\theta = (\theta_1, \theta_2, \dots, \theta_n)^T$ 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 $\partial LL / \partial \theta$ using high-school calculus
3. Solve the set of simultaneous equations

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

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

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

4. Check that you're at a maximum

A General MLE strategy

Suppose $\theta = (\theta_1, \theta_2, \dots, \theta_n)^T$ 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 $\partial LL / \partial \theta$ using high-school calculus
3. Solve the set of simultaneous equations

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

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

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

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

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 $x_1, x_2, \dots, x_R \sim (\text{i.i.d}) N(\mu, \sigma^2)$
- 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 $x_1, x_2, \dots, x_R \sim (\text{i.i.d}) N(\mu, \sigma^2)$
- 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 $E[\sigma_{\text{unbiased}}^2] = \sigma^2$

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

Unbiasedness discussion

- Which is best?

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

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

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

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

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

Are either of these unbiased?

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

Which is more useful?

MLE for m-dimensional Gaussian

- Suppose you have $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_R \sim (\text{i.i.d}) \mathcal{N}(\mu, \Sigma)$
- But you don't know μ or Σ
- MLE: For which $\theta = (\mu, \Sigma)$ is $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_R$ most likely?

$$\boldsymbol{\mu}^{mle} = \frac{1}{R} \sum_{k=1}^R \mathbf{x}_k$$

$$\boldsymbol{\Sigma}^{mle} = \frac{1}{R} \sum_{k=1}^R (\mathbf{x}_k - \boldsymbol{\mu}^{mle})(\mathbf{x}_k - \boldsymbol{\mu}^{mle})^T$$

MLE for m-dimensional Gaussian

- Suppose you have $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_R \sim (\text{i.i.d}) N(\mu, \Sigma)$
- But you don't know μ or Σ
- MLE: For which $\theta = (\mu, \Sigma)$ is $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_R$ most likely?

$$\boldsymbol{\mu}^{mle} = \frac{1}{R} \sum_{k=1}^R \mathbf{x}_k$$

$$\mu_i^{mle} = \frac{1}{R} \sum_{k=1}^R x_{ki}$$

$$\Sigma^{mle} = \frac{1}{R} \sum_{k=1}^R (\mathbf{x}_k - \boldsymbol{\mu}^{mle})(\mathbf{x}_k - \boldsymbol{\mu}^{mle})^T$$

Where $1 \leq i \leq m$

And x_{ki} is value of the i^{th} component of \mathbf{x}_k (the i^{th} attribute of the k^{th} record)

And μ_i^{mle} is the i^{th} component of $\boldsymbol{\mu}^{mle}$

MLE for m-dimensional Gaussian

- Suppose you have $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_R \sim (\text{i.i.d}) N(\mu, \Sigma)$
- But you don't know μ or Σ
- MLE: For which $\theta = (\mu, \Sigma)$ is $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_R$ most likely?

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

$$\Sigma^{mle} = \frac{1}{R} \sum_{k=1}^R (\mathbf{x}_k - \mu^{mle})(\mathbf{x}_k - \mu^{mle})^T$$

$$\sigma_{ij}^{mle} = \frac{1}{R} \sum_{k=1}^R (\mathbf{x}_{ki} - \mu_i^{mle})(\mathbf{x}_{kj} - \mu_j^{mle})$$

Where $1 \leq i \leq m, 1 \leq j \leq m$

And x_{ki} is value of the i^{th} component of \mathbf{x}_k (the i^{th} attribute of the k^{th} record)

And σ_{ij}^{mle} is the $(i,j)^{\text{th}}$ component of Σ^{mle}

MLE for m-dimensional Gaussian

- Suppose you have $\mathbf{x}_1, \mathbf{x}_2, \dots$
- But you don't know μ or Σ
- MLE: For which $\theta = (\mu, \Sigma)$ is \mathbf{x}

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

$$\Sigma^{mle} = \frac{1}{R} \sum_{k=1}^R (\mathbf{x}_k - \mu^{mle})(\mathbf{x}_k - \mu^{mle})^T$$

$$\Sigma^{\text{unbiased}} = \frac{\Sigma^{mle}}{1 - \frac{1}{R}} = \frac{1}{R-1} \sum_{k=1}^R (\mathbf{x}_k - \mu^{mle})(\mathbf{x}_k - \mu^{mle})^T$$

Q: How would you prove this?

A: Just plug through the MLE recipe.

Note how Σ^{mle} is forced to be symmetric non-negative definite

Note the unbiased case

How many datapoints would you need before the Gaussian has a chance of being non-degenerate?

Confidence intervals

We need to talk

We need to discuss how accurate we expect μ^{mle} and Σ^{mle} to be as a function of R

And we need to consider how to estimate these accuracies from data...

- Analytically *
- Non-parametrically (using randomization and bootstrapping) *

But we won't. Not yet.

*Will be discussed in future Andrew lectures...just before we need this technology.

Structural error

Actually, we need to talk about something else too..

What if we do all this analysis when the true distribution is in fact not Gaussian?

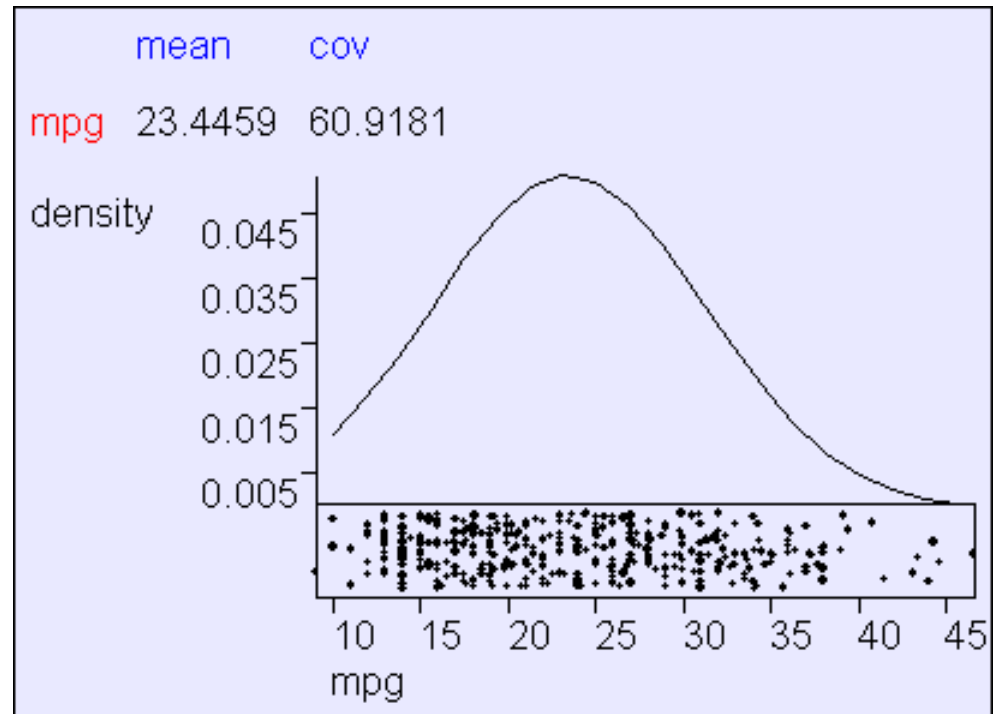
How can we tell? *

How can we survive? *

*Will be discussed in future Andrew lectures...just before we need this technology.

Gaussian MLE in action

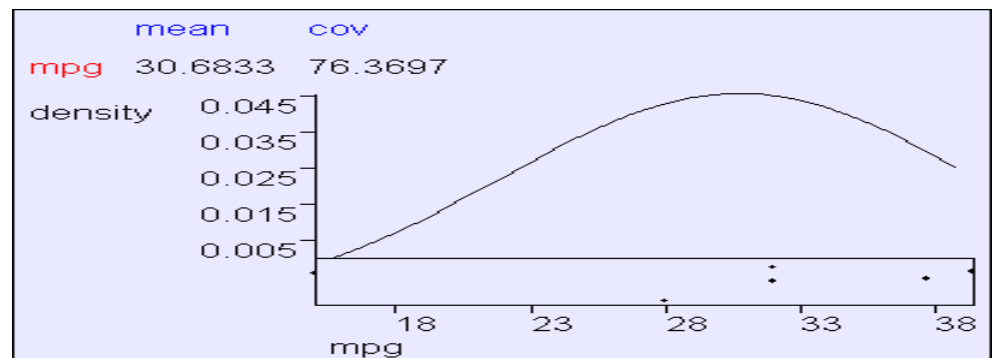
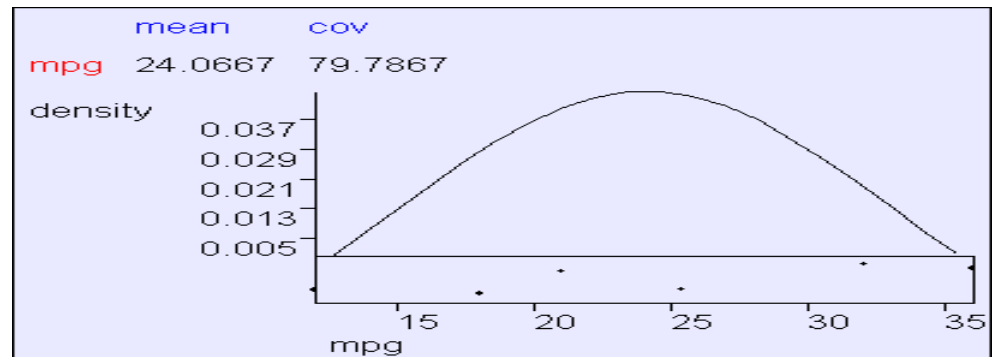
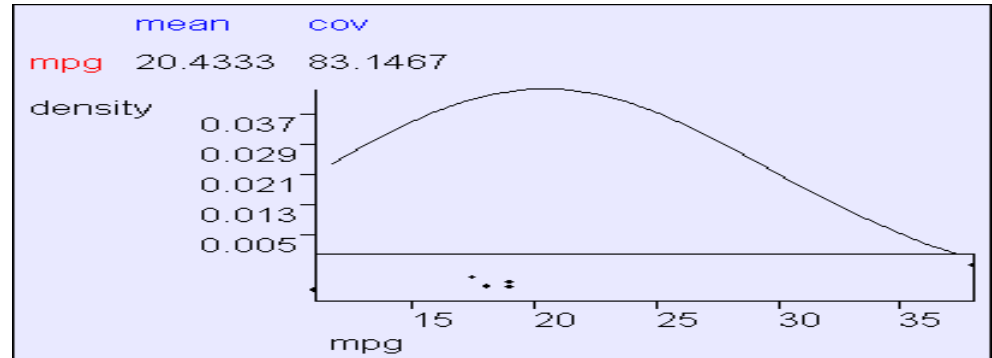
Using R=392 cars from the
“MPG” UCI dataset supplied
by Ross Quinlan



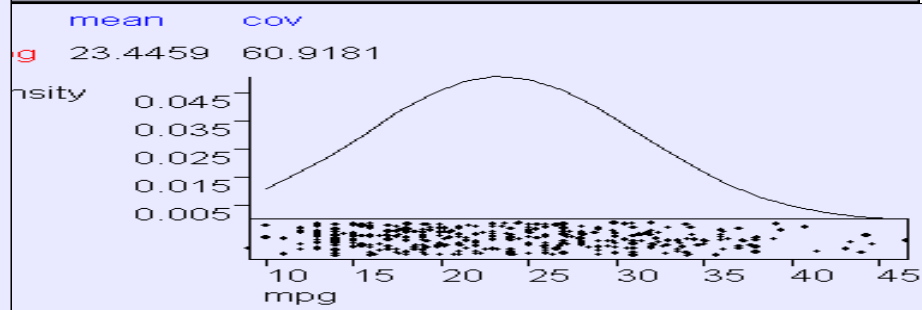
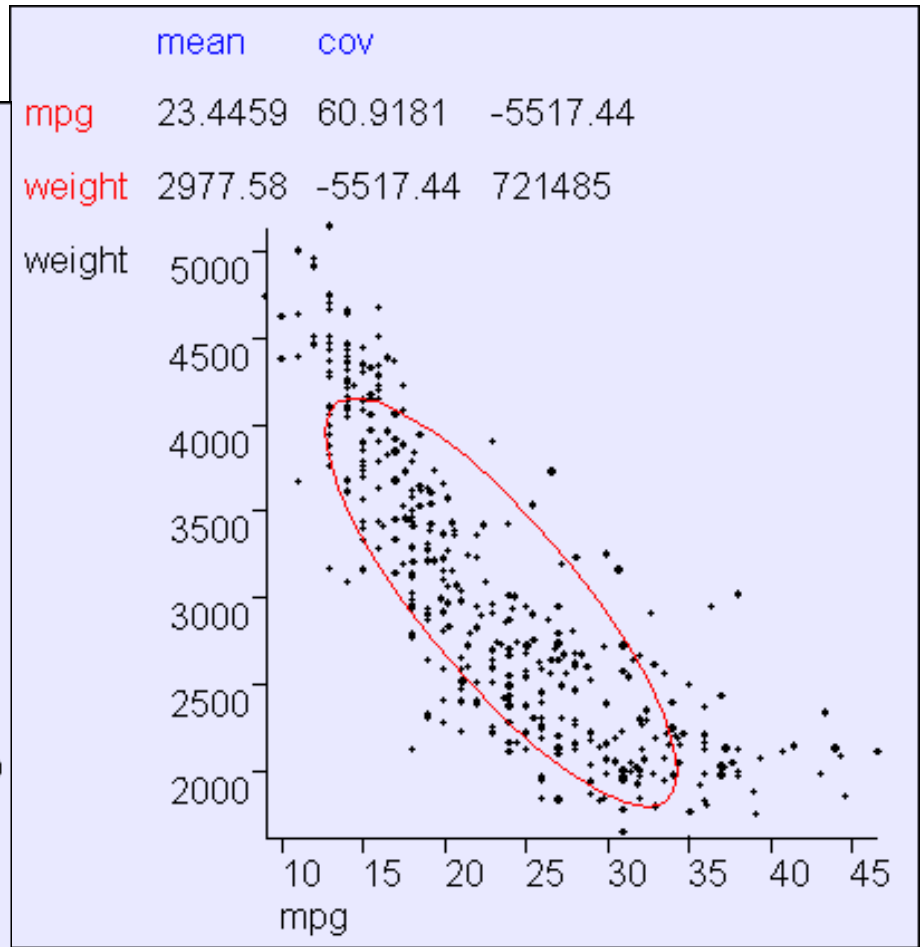
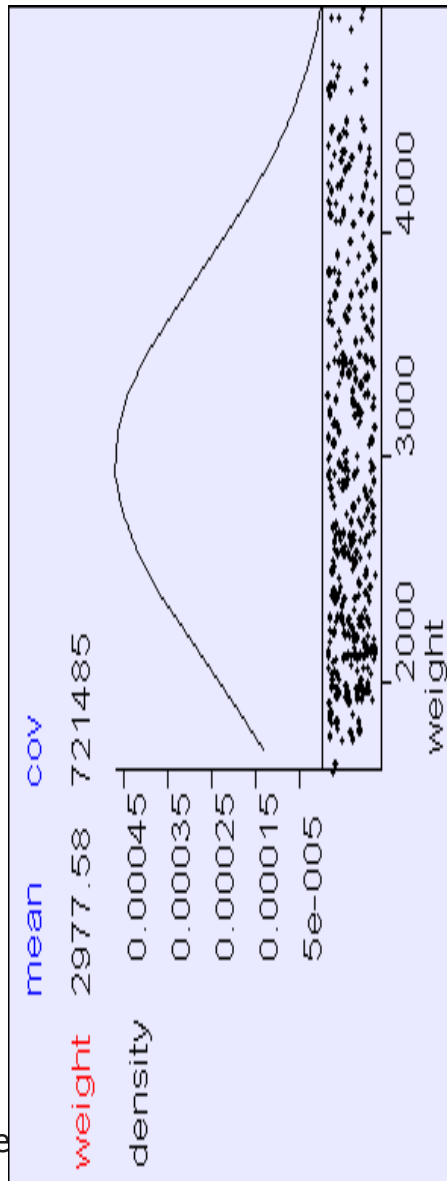
Data-starved Gaussian MLE

Using three subsets of MPG.

Each subset has 6
randomly-chosen cars.



Bivariate MLE in action



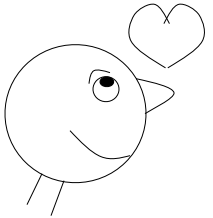
Multivariate MLE

	mean	cov						
mpg	23.4459	60.9181	-10.3529	-657.585	-233.858	-5517.44	9.11551	16.6915
cylinders	5.47194	-10.3529	2.9097	169.722	55.3482	1300.42	-2.37505	-2.17193
displacement	194.412	-657.585	169.722	10950.4	3614.03	82929.1	-156.994	-142.572
horsepower	104.469	-233.858	55.3482	3614.03	1481.57	28265.6	-73.187	-59.0364
weight	2977.58	-5517.44	1300.42	82929.1	28265.6	721485	-976.815	-967.228
acceleration	15.5413	9.11551	-2.37505	-156.994	-73.187	-976.815	7.61133	2.95046
modelyear	75.9796	16.6915	-2.17193	-142.572	-59.0364	-967.228	2.95046	13.5699

Covariance matrices are not exciting to look at

Being Bayesian: MAP estimates for Gaussians

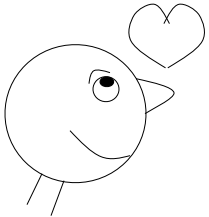
- Suppose you have $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_R \sim (\text{i.i.d}) \mathcal{N}(\mu, \Sigma)$
- But you don't know μ or Σ
- MAP: Which (μ, Σ) maximizes $p(\mu, \Sigma | \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_R)$?



Step 1: Put a prior on (μ, Σ)

Being Bayesian: MAP estimates for Gaussians

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



Step 1: Put a prior on (μ, Σ)

Step 1a: Put a prior on Σ

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

This thing is called the Inverse-Wishart distribution.

A PDF over SPD matrices!

Doing Bayesian: MAP estimates for Gaussians

- v_0 small: "I am not sure about my guess of Σ_0 "

- v_0 large: "I'm pretty sure about my guess of Σ_0 "

Σ_0 : (Roughly) my best guess of Σ

$$E[\Sigma] = \Sigma_0$$

Step 1: Put a prior on (μ, Σ)

Step 1a: Put a prior on Σ

$$(v_0 - m - 1) \Sigma \sim \text{IW}(v_0, (v_0 - m - 1) \Sigma_0)$$

This thing is called the Inverse-Wishart distribution.

A PDF over SPD matrices!

Being Bayesian: MAP estimates for Gaussians

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

Step 1: Put a prior on (μ, Σ)

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 $\mu | \Sigma$

$$\mu | \Sigma \sim \mathcal{N}(\mu_0, \Sigma / \kappa_0)$$

Together, " Σ " and " $\mu | \Sigma$ " define a joint distribution on (μ, Σ)

Being Bayesian: MAP estimates for Gaussians

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

μ_0 : My best guess of μ (μ, Σ)
 $E[\mu] = \mu_0$

κ_0 small: "I am not sure
 about my guess of μ_0 "

κ_0 large: "I'm pretty sure
 about my guess of μ_0 "

$(v_0 - m - 1)\Sigma \sim \dots, v_0, (v_0 - m - 1)\Sigma$

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

$$\mu \mid \Sigma \sim N(\mu_0, \Sigma / \kappa_0)$$

Together, " Σ " and
 " $\mu \mid \Sigma$ " define a
 joint distribution
 on (μ, Σ)

Notice how we are forced to express our
 ignorance of μ proportionally to Σ

Being Bayesian: MAP estimates for Gaussians

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

Step 1: Put a prior on (μ, Σ)

Why do we use this form of prior?

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 $\mu | \Sigma$

$$\mu | \Sigma \sim \mathcal{N}(\mu_0, \Sigma / \kappa_0)$$

Being Bayesian: MAP estimates for Gaussians

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

Step 1: Put a prior on (μ, Σ)

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 $\mu | \Sigma$

$$\mu | \Sigma \sim \mathcal{N}(\mu_0, \Sigma / \kappa_0)$$

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

Step 1: Prior: $(\nu_0 - m - 1) \Sigma \sim \text{IW}(\nu_0, (\nu_0 - m - 1) \Sigma_0)$, $\mu \mid \Sigma \sim \mathcal{N}(\mu_0, \Sigma / \kappa_0)$

Step 2:

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

$$\mu_R = \frac{\kappa_0 \mu_0 + R \bar{\mathbf{x}}}{\kappa_0 + R}$$

$$\nu_R = \nu_0 + R$$

$$\kappa_R = \kappa_0 + R$$

$$(\nu_R + m - 1) \Sigma_R = (\nu_0 + m - 1) \Sigma_0 + \sum_{k=1}^R (\mathbf{x}_k - \bar{\mathbf{x}})(\mathbf{x}_k - \bar{\mathbf{x}})^T + \frac{(\bar{\mathbf{x}} - \mu_0)(\bar{\mathbf{x}} - \mu_0)^T}{1/\kappa_0 + 1/R}$$

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

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

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

Being Bayesian

- Suppose you have

- MAP: Which (μ, Σ)

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

Step 2:

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

$$\mu_R = \frac{\kappa_0 \mu_0 + R \bar{\mathbf{x}}}{\kappa_0 + R}$$

$$\nu_R = \nu_0 + R$$

$$\kappa_R = \kappa_0 + R$$

• 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 μ is a student-t

• One point of view: it's pretty academic if $R > 30$

$$(\nu_R + m - 1) \Sigma_R = (\nu_0 + m - 1) \Sigma_0 + \sum_{k=1}^R (\mathbf{x}_k - \bar{\mathbf{x}})(\mathbf{x}_k - \bar{\mathbf{x}})^T + \frac{(\bar{\mathbf{x}} - \mu_0)(\bar{\mathbf{x}} - \mu_0)^T}{1/\kappa_0 + 1/R}$$

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

$$\mu \mid \Sigma \sim N(\mu_R, \Sigma / \kappa_R)$$

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

Where we're at

			Categorical inputs only	Real-valued inputs only	Mixed Real / Cat okay
Inputs	<div>Classifier</div> <div>Predict category</div>	Joint BC Naïve BC			Dec Tree
Inputs	<div>Density Estimator</div> <div>Probability</div>	Joint DE Naïve DE		Gauss DE	
Inputs	<div>Regressor</div> <div>Predict real no.</div>				

What you should know

- The Recipe for MLE
- What do we sometimes prefer MLE to MAP?
- Understand MLE estimation of Gaussian parameters
- Understand “biased estimator” versus “unbiased estimator”
- Appreciate the outline behind Bayesian estimation of Gaussian parameters

Useful exercise

- We'd already done some MLE in this class without even telling you!
- Suppose categorical arity- n inputs $x_1, x_2, \dots, x_R \sim (\text{i.i.d.})$ from a multinomial

$$M(p_1, p_2, \dots, p_n)$$

where

$$P(x_k=j|\mathbf{p})=p_j$$

- What is the MLE $\mathbf{p}=(p_1, p_2, \dots, p_n)$?