

CS 5135/6035 Learning Probabilistic Models

Lecture 8: MLE, Gradient Descent, Multivariate Gaussian

Gowtham Atluri

September 17, 2018

Gowtham Atluri CS 5135/6035 Learning Probabilistic Models September 17, 2018 1 / 35

Reading Material

- Properties of MLE

<https://ocw.mit.edu/courses/mathematics/18-443-statistics-for-applications-fall-2006/lecture-notes/lecture3.pdf>

- Singer, Advanced Optimization, Lecture 9

https://people.seas.harvard.edu/~yaron/AM221-S16/lecture_notes/AM221_lecture9.pdf

- Jordan, Chapter 13. The Multivariate Gaussian

<https://people.eecs.berkeley.edu/~jordan/courses/260-spring10/other-readings/chapter13.pdf>

Gowtham Atluri CS 5135/6035 Learning Probabilistic Models September 17, 2018 2 / 35

Maximum Likelihood Estimation

- Review
- Properties of Estimators
 - Consistency
 - Bias
 - Variance

Gowtham Atluri CS 5135/6035 Learning Probabilistic Models September 17, 2018 3 / 35

Maximum Likelihood Estimation - Recap

- I.I.D assumption
 - x_1, x_2, \dots, x_n are i.i.d.
 - x_i 's are *independently* sampled
 - Ever x_i is drawn from the *same* probability distribution
- Likelihood

$$p(x_1, \dots, x_n | \theta) = p(x_1 | \theta) p(x_2 | \theta) \dots p(x_n | \theta) = \prod_{i=1}^n f(x_i | \theta) = L(\theta | x)$$

- Log-likelihood

$$\ell(\theta) = \log L(\theta)$$

- Maximization

$$\hat{\theta}_{MLE} = \operatorname{argmax}_{\theta} p(D | \theta) = \operatorname{argmax}_{\theta} \ell(\theta)$$

- Examples
 - Bernoulli (Discrete)
 - Gaussian (Continuous)

Gowtham Atluri CS 5135/6035 Learning Probabilistic Models September 17, 2018 4 / 35

Properties of Estimators - Recap

Consistency

An estimator is consistent if the estimate $\hat{\theta}$ it constructs is guaranteed to converge to the true parameter value θ as the quantity of data to which it is applied increases.

Bias

The bias of an estimator η is defined as the deviation of the expectation of the estimate from the true value: $E[\hat{\theta}_{\eta}]$

When the sampling of data is viewed as a stochastic process, then the estimated parameter $\hat{\theta}_{\eta}$ can be viewed as a continuous random variable.

When $E[\hat{\theta}_{\eta}] = \theta$ we say the estimator is unbiased.}

Variance (and efficiency)

$\operatorname{Var}[\hat{\theta}_{\eta}]$

All else being equal, an estimator with smaller variance is preferable to one with greater variance.

Gowtham Atluri CS 5135/6035 Learning Probabilistic Models September 17, 2018 5 / 35

MLE for univariate Gaussian variable

The temperatures, in Celsius, in Minneapolis during the first week of March 2018 are observed as $(-2.5, -9.9, -12.1, -8.9, -6.0, -4.8, 2.4)$
What is the distribution from which this data was generated (assuming it was Gaussian)?

$$\hat{\mu} = \frac{\sum_{i=1}^7 x_i}{7} = -5.97$$

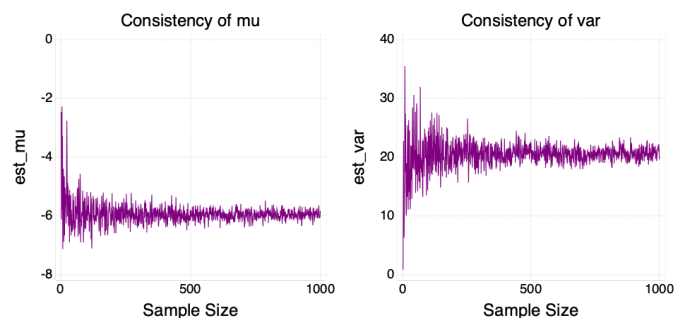
$$\hat{\sigma}^2 = \frac{\sum_{i=1}^7 (x_i - \mu)^2}{7} = 20.72$$

Gowtham Atluri CS 5135/6035 Learning Probabilistic Models September 17, 2018 6 / 35

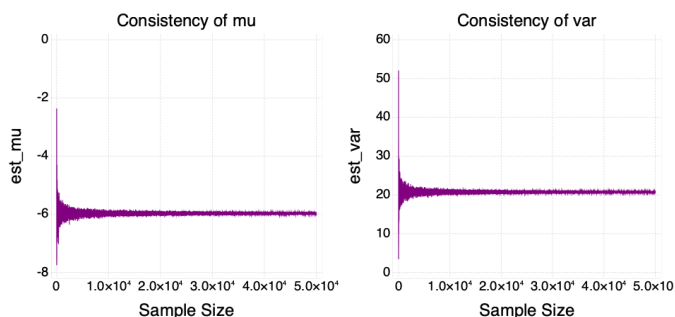
Estimating Consistency (Julia)

```
d = Normal(-5.97,sqrt(20.72));
sample_size = collect(2:1000);
est_mean = zeros(length(sample_size));
est_var = zeros(length(sample_size));
for i=1:length(sample_size)
    sample = rand(d,sample_size[i]);
    est_mean[i] = sum(sample)/length(sample);
    est_var[i] = sum((sample.-est_mean[i]).^2)/length(sample);
end
myplot1 = Gadfly.plot(x=sample_size,y=est_mean,Geom.line,
    Guide.xlabel("Sample Size"), Guide.ylabel("est_mu"),
    Guide.title("Consistency of mu"),white_panel);
myplot2 = Gadfly.plot(x=sample_size,y=est_var,Geom.line,
    Guide.xlabel("Sample Size"), Guide.ylabel("est_var"),
    Guide.title("Consistency of var"),white_panel);
final_plot = hstack(myplot1,myplot2);
draw(PNG("./figs/univ_normal_consistency.png", 10inch, 5inch), final_plot)
```

Estimating Consistency (Julia)



Estimating Consistency (Julia)



Estimating Bias ($\hat{\mu}$)

Estimator of mean is unbiased if $E[\hat{\theta}_{MLE}] = \theta$

$$\begin{aligned}\hat{\mu} &= \frac{\sum_{i=1}^n x_i}{n} \\ E(\hat{\mu}) &= E\left[\frac{1}{n} \sum_{i=1}^n x_i\right] \\ &= \frac{1}{n} E\left[\sum_{i=1}^n x_i\right] \\ &= \frac{1}{n} \cdot n \cdot E[x] \\ &= \mu \quad (\text{as } E(x) = \mu)\end{aligned}$$

The mean estimator is unbiased.

Estimating Bias ($\hat{\sigma}^2$)

Estimator of variance is unbiased if $E[\hat{\sigma}_{MLE}^2] = \sigma^2$

$$\begin{aligned}\hat{\sigma}^2 &= \frac{\sum_{i=1}^n (x_i - \mu)^2}{n} \\ E(\hat{\sigma}^2) &= E\left(\frac{\sum_{i=1}^n (x_i - \mu)^2}{n}\right) \\ &= \dots \\ &= \frac{n-1}{n} \sigma^2\end{aligned}$$

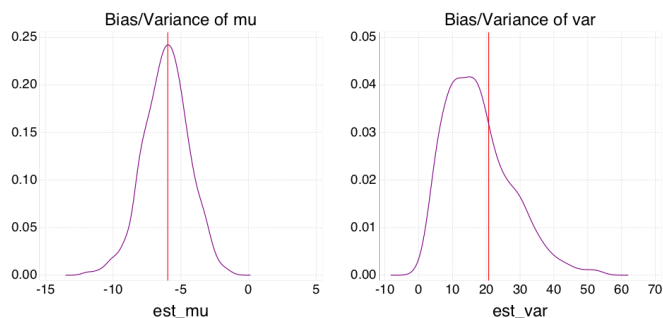
The variance estimator is biased.

Estimating Bias and Variance (Julia)

```
d = Normal(-5.97,sqrt(20.72));
est_mean = zeros(1000);
est_var = zeros(1000);
for i=1:1000
    sample = rand(d,7);
    est_mean[i] = sum(sample)/length(sample);
    est_var[i] = sum((sample.-est_mean[i]).^2)/length(sample);
end
myplot1 = Gadfly.plot(x=est_mean,Geom.density, Guide.xlabel("est_mu",
    xintercept=[-5.97], Geom.vline(color=colorant"red"),
    Guide.title("Bias/Variance of mu"),white_panel);
myplot2 = Gadfly.plot(x=est_var,Geom.density, Guide.xlabel("est_var",
    xintercept=[20.72], Geom.vline(color=colorant"red"),
    Guide.title("Bias/Variance of var"),white_panel);
final_plot = hstack(myplot1,myplot2);
draw(PNG("./figs/univ_normal_bias_var_7.png", 10inch, 5inch), final_plot)
```

Estimating Bias and Variance (Julia)

$n = 7$



Gowham Atluri

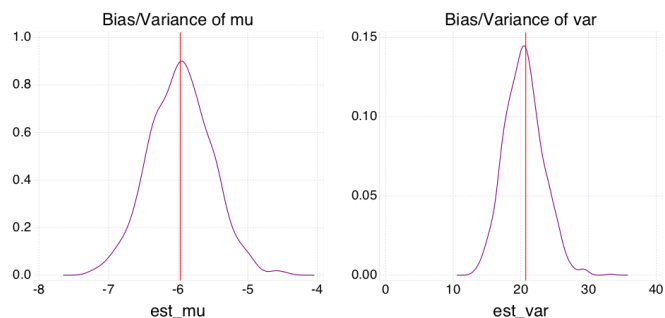
CS 5135/6035 Learning Probabilistic Models

September 17, 2018

13 / 35

Estimating Bias and Variance (Julia)

$n = 100$



Gowham Atluri

CS 5135/6035 Learning Probabilistic Models

September 17, 2018

14 / 35

Gradient Descent Approach for MLE

- Overview
- General algorithm
- MLE for Gamma Distr.
- Julia code
- Limitations of Gradient Descent

Gowham Atluri

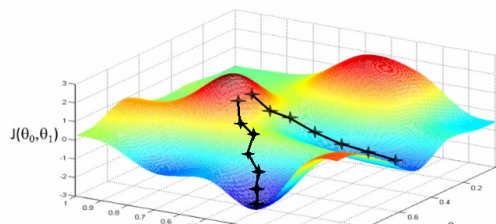
CS 5135/6035 Learning Probabilistic Models

September 17, 2018

15 / 35

Gradient Descent for MLE: Approach - II

- When approach I is not possible (particularly when the model involves many parameters and its PDF is highly non-linear), use gradient descent approach.
 - Use negative log-likelihood (also referred to as a cost function)
 - Randomly initialize and then incrementally update our weights by calculating the slope of our objective function
 - When applying the cost function, we want to continue updating our weights until the slope of the gradient gets as close to zero as possible.



Gowham Atluri

CS 5135/6035 Learning Probabilistic Models

September 17, 2018

16 / 35

Gradient Descent

- A numerical optimization technique used to find the parameter vector \mathbf{w} that minimizes an objective function $E(\mathbf{w})$.

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} E(\mathbf{w})$$

- An iterative approach to estimate \mathbf{w}
 - Starts with an initial estimate \mathbf{w}_1 (often a random vector)
 - Steepest descent from this point is to follow the negative gradient $-\nabla E$

$$\nabla E = \left[\frac{dE}{dw_1} \quad \frac{dE}{dw_2} \quad \dots \quad \frac{dE}{dw_n} \right]$$

- Next estimate \mathbf{w}_2 is estimated as $\mathbf{w}_2 \leftarrow \mathbf{w}_1 - \lambda \nabla E|_{\mathbf{w}_1}$
- Generally $\mathbf{w}_i \leftarrow \mathbf{w}_{i-1} - \lambda \nabla E|_{\mathbf{w}_{i-1}}$
 - λ is the learning rate
- Stops after a given maxIter or when estimate \mathbf{w} converges

Gowham Atluri

CS 5135/6035 Learning Probabilistic Models

September 17, 2018

17 / 35

Gradient Descent: a general algorithm

- Step 1: Pick initial value \mathbf{w}_1
- Step 2: $\text{maxIter} = 10000$
- Step 3: **for** $i = 2 : \text{maxIter}$
- Step 4: $\mathbf{w}_i \leftarrow \mathbf{w}_{i-1} - \lambda \nabla E|_{\mathbf{w}_{i-1}}$
- Step 5: **if** $|\mathbf{w}_i - \mathbf{w}_{i-1}| < \epsilon$ **terminate; end**
- Step 6: **end for**

Gowham Atluri

CS 5135/6035 Learning Probabilistic Models

September 17, 2018

18 / 35

MLE for Gamma distribution

Probability density function of Gamma distribution is

$$\frac{1}{\Gamma(\alpha)\beta^\alpha} x^{\alpha-1} e^{-x/\beta}$$

where $\Gamma(\alpha)$ is the gamma function and (α, β) are parameters that take positive values.

Likelihood function

$$L(\theta|x) = \frac{1}{\Gamma(\alpha)^n \beta^{n\alpha}} \left(\prod_i x_i^{\alpha-1} \right) e^{-\sum_i x_i/\beta}$$

Log-Likelihood function

$$\ell(\theta) = -n \log \Gamma(\alpha) - n\alpha \log \beta + (\alpha - 1) \sum_i \log x_i - \frac{\sum_i x_i}{\beta}$$

Negative Log-Likelihood function

$$-\ell(\theta) = n \log \Gamma(\alpha) + n\alpha \log \beta - (\alpha - 1) \sum_i \log x_i + \frac{\sum_i x_i}{\beta}$$

Gowtham Atluri

CS 5135/6035 Learning Probabilistic Models

September 17, 2018

19 / 35

MLE for Gamma distribution

Negative Log-Likelihood function

$$-\ell(\theta) = n \log \Gamma(\alpha) + n\alpha \log \beta - (\alpha - 1) \sum_i \log x_i + \frac{\sum_i x_i}{\beta}$$

Computing partial derivatives:

$$\frac{\partial \ell}{\partial \alpha} = n \frac{\partial}{\partial \alpha} \log \Gamma(\alpha) + n \log \beta - \sum_i \log x_i$$

$$\frac{\partial \ell}{\partial \beta} = n \frac{\alpha}{\beta} - \frac{\sum_i x_i}{\beta^2}$$

Gradient Descent update rules:

$$\alpha \leftarrow \alpha - \gamma \frac{\partial \ell}{\partial \alpha} \quad \beta \leftarrow \beta - \gamma \frac{\partial \ell}{\partial \beta}$$

where γ is the learning rate.

Gowtham Atluri

CS 5135/6035 Learning Probabilistic Models

September 17, 2018

20 / 35

MLE using Gradient Descent (Julia)

```
# generating 10 samples from a Gamma distribution
d = Gamma(5,5);
sample = rand(d,1000)
```

```
## 1000-element Array{Float64,1}:
##  22.5362
##  25.1435
##  45.6914
##  26.7435
##  37.8724
##  26.1113
##  21.9101
##  18.8474
##  16.8973
##  23.8151
##
##  18.313
##  16.7582
##  6.79626
```

Gowtham Atluri

CS 5135/6035 Learning Probabilistic Models

September 17, 2018

21 / 35

MLE using Gradient Descent (Julia)

```
function dl_by_da(sample,a,b)
    n = length(sample);
    result = n*digamma(a) + n*log(b) - sum(log.(sample));
    return result;
end
```

dl_by_da (generic function with 1 method)

```
function dl_by_db(sample,a,b)
    n = length(sample);
    result = (n*a/b) - (sum(sample)/(b^2));
    return result;
end
```

dl_by_db (generic function with 1 method)

Gowtham Atluri

CS 5135/6035 Learning Probabilistic Models

September 17, 2018

22 / 35

MLE using Gradient Descent (Julia)

```
function gradient_descent_gamma(sample)
    n = length(sample);
    max_itr = 1000; # maximum num. iterations
    gm = 0.01; # rate of learning
    a = rand()*10; # random initialization
    b = rand()*10;

    for i=1:max_itr
        a_new = a - gm*dl_by_da(sample,a,b);
        b_new = b - gm*dl_by_db(sample,a,b);
        if(a_new<0) a_new = rand()*10; end;
        if(b_new<0) b_new = rand()*10; end;
        if(abs(a_new-a)<0.0001 && abs(b_new-b)<0.0001) break; end;
        a = a_new; b = b_new;
    end
    return a,b;
end
```

```
## gradient_descent_gamma (generic function with 1 method)
```

Gowtham Atluri

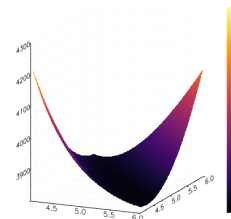
CS 5135/6035 Learning Probabilistic Models

September 17, 2018

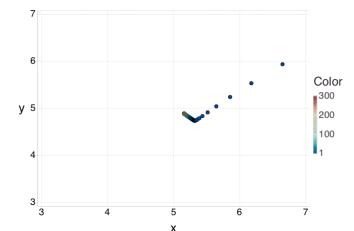
23 / 35

MLE using Gradient Descent (Julia)

Shape of $-\ell(\theta)$ around (5,5)



Optimization path



Gowtham Atluri

CS 5135/6035 Learning Probabilistic Models

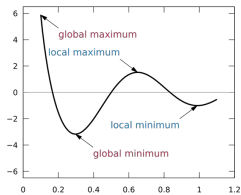
September 17, 2018

24 / 35

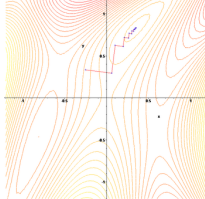
Gradient Descent: limitation

- Can converge to a local minimum
 - can result in a different value in different runs
- Tends to be slow when it is close to the minimum
- In poorly conditioned convex problems, 'zigzags' when gradients point nearly orthogonally to the shortest direction

Convergence to local minimum



Zigzag gradients



Multivariate Gaussian

- Functional form
- Covariance matrix
- Isocontours
- Multivariate Gaussian as a product of univariate distributions
- Properties of Multivariate Gaussian

Multivariate Gaussian

Univariate Gaussian

$$p(x|\mu, \sigma^2) = \mathcal{N}(x|\mu, \sigma^2) \equiv \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

Multivariate Gaussian

$$p(\mathbf{x}|\mu, \Sigma) = \mathcal{N}(\mathbf{x}|\mu, \Sigma) \equiv \frac{1}{\sqrt{\det(2\pi\Sigma)}} e^{-\frac{1}{2}(\mathbf{x}-\mu)^T \Sigma^{-1}(\mathbf{x}-\mu)}$$

Univariate

- exponent $-\frac{1}{2\sigma^2}(x-\mu)^2$
 - quadratic in x
 - negative sign
- coefficient in front $\frac{1}{\sqrt{2\pi\sigma^2}}$
 - normalization

Multivariate

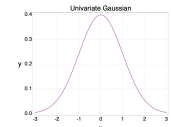
- $-\frac{1}{2}(\mathbf{x}-\mu)^T \Sigma^{-1}(\mathbf{x}-\mu)$
 - quadratic in \mathbf{x}
 - negative sign
- $\frac{1}{\sqrt{\det(2\pi\Sigma)}}$
 - normalization

Multivariate Gaussian

```
rand(Normal(0,1),4)
```

```
## 4-element Array{Float64,1}:
##  0.654719
## -1.36141
##  0.830842
##  0.192311
```

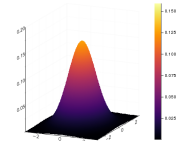
Univariate Gaussian



```
rand(MvNormal([0,0],eye(2)), 4)'
```

```
## 4x2 Array{Float64,2}:
## -0.792111 -0.465018
## -1.02169  0.512884
## -1.10957  0.821939
##  0.271616 -1.05763
```

Multivariate Gaussian



Multivariate Gaussian

$$p(\mathbf{x}|\mu, \Sigma) = \mathcal{N}(\mathbf{x}|\mu, \Sigma) \equiv \frac{1}{\sqrt{\det(2\pi\Sigma)}} e^{-\frac{1}{2}(\mathbf{x}-\mu)^T \Sigma^{-1}(\mathbf{x}-\mu)}$$

where μ is the mean vector of the distribution

and Σ is the covariance matrix.

Covariance Matrix:

- For two random variables x, y ,

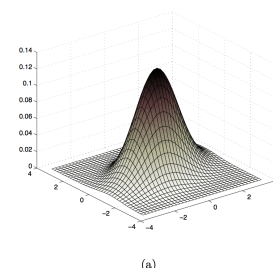
$$\text{Cov}[x, y] = E[(x - E(x))(y - E(y))] = E[xy] - E[x]E[y]$$
 - $\Sigma_{ij} = \text{Cov}(x_i, x_j)$
- $\Sigma \in \mathbf{S}^n_{++}$ (Positive Definite)

$$\mathbf{S}^n_{++} = \{A \in \mathbf{R}^{n \times n} : A = A^T \text{ and } x^T A x > 0 \ \forall x \in \mathbf{R}^n, \text{ such that } x \neq 0\}$$
 - If all eigenvalues are positive, then the matrix is positive definite

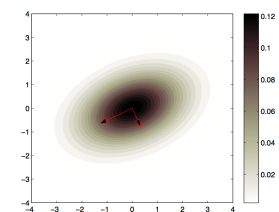
Isocontours

For a function $f: \mathbf{R}^2 \rightarrow \mathbf{R}$, an isocontour is a set of the form

$$x \in \mathbf{R}^2 : f(x) = c, \text{ for some } c \in \mathbf{R}$$



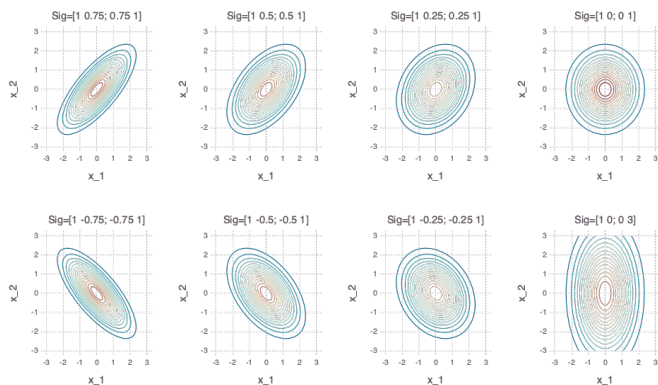
(a)



(b)

Mean (0, 0), Covariance $[1, 0.5; 0.5, 1.75]$

Multivariate Gaussian: Covariance matrix - Geometric view



Gowham Atluri

CS 5135/6035 Learning Probabilistic Models

September 17, 2018

31 / 35

Multivariate Gaussian - Geometric view

- Every real symmetric matrix $D \times D$ has an eigen-decomposition

$$\Sigma = E\Lambda E^T$$

where $E^T E = I$ and $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_D)$

- In the case of covariance matrix, all eigenvalues λ_i are positive.
- One can then use

$$y = \Lambda^{-\frac{1}{2}} E^T (x - \mu)$$

so that

$$(x - \mu)^T \Sigma^{-1} (x - \mu) = (x - \mu)^T E \Lambda^{-1} E^T (x - \mu) = y^T y$$

- The multivariate Gaussian reduces to a product of D univariate zero-mean unit variance Gaussians.

$$X \sim \mathcal{N}(\mu, \Sigma) \iff X \sim \mu + E\mathcal{N}(0, \Lambda^{1/2}) \iff X \sim \mu + E\Lambda^{1/2}\mathcal{N}(0, I)$$

- We can view multivariate Gaussian as a shifted, scaled, and rotated version of 'standard' Gaussian in which the center is given by the mean, rotation by the eigen vectors and scaling by sqroot of

Gowham Atluri

CS 5135/6035 Learning Probabilistic Models

September 17, 2018

32 / 35

Product of Gaussians

Product of two Gaussians is another Gaussian with a multiplicative factor

$$\mathcal{N}(x|\mu_1, \Sigma_1)\mathcal{N}(x|\mu_2, \Sigma_2) = \mathcal{N}(x|\mu, \Sigma) \frac{\exp(-\frac{1}{2}(\mu_1 - \mu_2)^T S^{-1}(\mu_1 - \mu_2))}{\sqrt{\det(2\pi S)}}$$

where $S \equiv \Sigma_1 + \Sigma_2$ and the mean and covariance are given by

$$\mu = \Sigma_1 S^{-1} \mu_2 + \Sigma_2 S^{-1} \mu_1 \quad \Sigma = \Sigma_1 S^{-1} \Sigma_2$$

Gowham Atluri

CS 5135/6035 Learning Probabilistic Models

September 17, 2018

33 / 35

Linear Transform of a Gaussian

Let y be linearly related to x through

$$y = Mx + \eta$$

where $x \perp \eta$, $\eta \sim \mathcal{N}(\mu, \Sigma)$ and $x \sim \mathcal{N}(\mu_x, \Sigma_x)$

Then marginal $p(y) = \int_x p(y|x)p(x)$ is a Gaussian

$$p(y) = \mathcal{N}(y|M\mu_x + \mu, M\Sigma_x M^T + \Sigma)$$

Gowham Atluri

CS 5135/6035 Learning Probabilistic Models

September 17, 2018

34 / 35

Partitioned Gaussian

Consider a distribution $\mathcal{N}(z|\mu, \Sigma)$ defined jointly over two vectors x and y of potentially different dimensions,

$$z = \begin{pmatrix} x \\ y \end{pmatrix}$$

with corresponding mean and partitioned covariance

$$\mu = \begin{pmatrix} \mu_x \\ \mu_y \end{pmatrix}, \quad \Sigma = \begin{pmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{pmatrix}$$

where $\Sigma_{yx} \equiv \Sigma_{xy}^T$.

The marginal distribution is given by

$$p(x) = \mathcal{N}(x|\mu_x, \Sigma_{xx})$$

and conditional

$$p(y|x) = \mathcal{N}(y|\mu_x + \Sigma_{xy}\Sigma_{yy}^{-1}(y - \mu_y), \Sigma_{xx} - \Sigma_{xy}\Sigma_{yy}^{-1}\Sigma_{yx})$$

Gowham Atluri

CS 5135/6035 Learning Probabilistic Models

September 17, 2018

35 / 35