# Probabilistic models for neural data:
# From single neurons to population dynamics

## NEUROBIO 316QC

Jan Drugowitsch

jan_drugowitsch@hms.harvard.edu

**Session 2**: Gaussians & Linear models

# Today

Q&A about previous session

Discuss assignment (~20min)

Gaussian distribution, linear models & model comparison (remaining time)

**Note**
no course on March 16
last session on March 30

# Overview

**Univariate and multivariate Gaussian distributions**

Means and variances of linearly transformed random variables (general)

Univariate standard Gaussian and linear transformations thereof

Multivariate Gaussians through linear constructions

Probabilistic operations: marginalization and conditioning

Maximum likelihood parameter estimation

Priors as added observations / regularizers

**Linear models**

Maximum likelihood and least squares estimates

What is linear in linear models?

Priors and Bayesian inference

**Bayesian model comparison**

Trading off goodness-of-fit with model complexity

Adjusting model complexity through hyperpriors

Evidence approximation

# Overview

**Univariate and multivariate Gaussian distributions**

　　Means and variances of linearly transformed random variables (general)

　　Univariate standard Gaussian and linear transformations thereof

　　Multivariate Gaussians through linear constructions

　　Probabilistic operations: marginalization and conditioning

　　Maximum likelihood parameter estimation

　　Priors as added observations / regularizers

**Linear models**

　　Maximum likelihood and least squares estimates

　　What is linear in linear models?

　　Priors and Bayesian inference

**Bayesian model comparison**

　　Trading off goodness-of-fit with model complexity

　　Adjusting model complexity through hyperpriors

　　Evidence approximation

# Mean and variance of linearly transformed RVs

**Linear transformation** of arbitrary single random variable $z$

Scaling

$x = az$

$$E[x] = \int x\, p(x)\, \mathrm{d}x = \int \overbrace{az}^{x}\, \overbrace{p(z)}^{p(x)} \overbrace{\left|\frac{\mathrm{d}z}{\mathrm{d}x}\right| \frac{\mathrm{d}x}{\mathrm{d}z}}^{\mathrm{d}x}\, \mathrm{d}z = a \int z\, p(z) \frac{1}{a} a\, \mathrm{d}z = a\, E[z]$$

$$\mathrm{var}[x] = a^2 \mathrm{var}[z]$$

Shifting & scaling

$x = az + b$

$$E[x] = a\, E[z] + b$$

$$\mathrm{var}[x] = a^2 \mathrm{var}[z]$$

**Linear transformation** of arbitrary random variable pair, $z_1$ and $z_2$
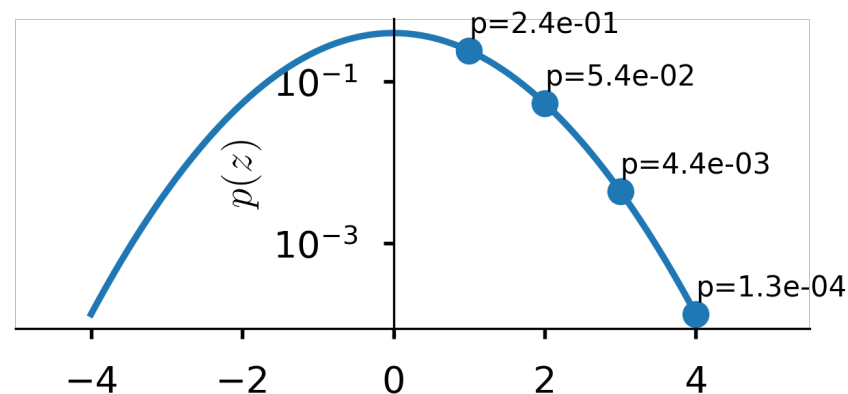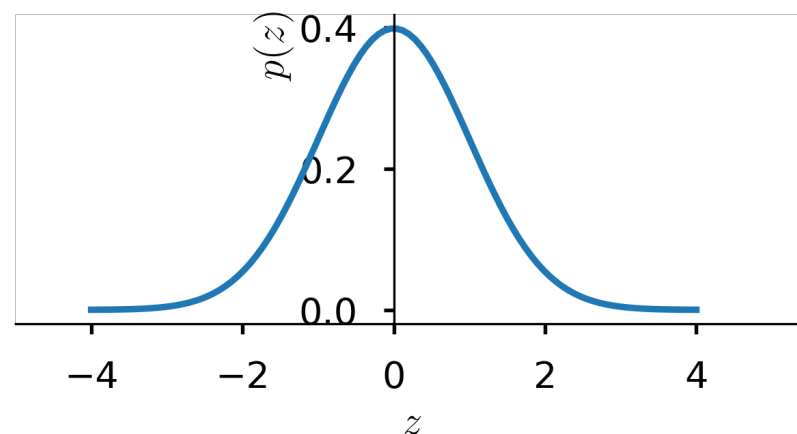
$$x = a_1 z_1 + a_2 z_2 + b$$

$$E[x] = a_1 E[z_1] + a_2 E[z_2] + b$$

$$var[x] = a_1^2\, var[z_1] + a_2^2\, var[z_2] + 2 a_1 a_2\, cov[z_1, z_2]$$

# The univariate standard Gaussian

**Standard Gaussian**: zero mean, unit variance

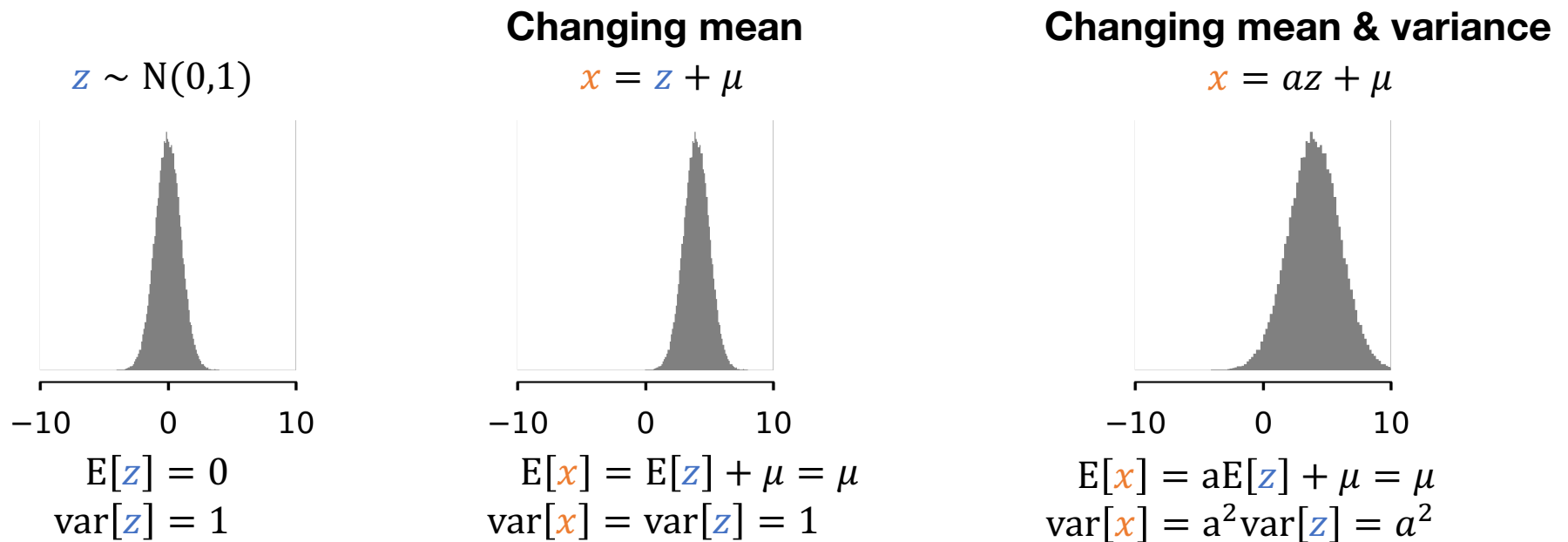mean  variance

$$p(z) = \mathrm{N}(z|0,1) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}z^2} \propto e^{-\frac{1}{2}z^2}$$

$$\log p(z) = -\frac{1}{2}z^2 - \frac{1}{2}\log 2\pi = -\frac{1}{2}z^2 + \mathrm{const.}$$



Rapidly dropping probabilities in tails:
<span style="color:red">considers outliers very unlikely</span>
(sensitivity to outliers)

# Linear transformations of univariate Gaussians

**Changing mean**

**Changing mean & variance**

$z \sim N(0,1)$

$x = z + \mu$

$x = az + \mu$



$$E[z] = 0$$
$$\text{var}[z] = 1$$

$$E[x] = E[z] + \mu = \mu$$
$$\text{var}[x] = \text{var}[z] = 1$$

$$E[x] = aE[z] + \mu = \mu$$
$$\text{var}[x] = a^2\text{var}[z] = a^2$$

## Linear transformation of Gaussian remains Gaussian

$$x = f(z) = az + \mu$$

$$f'(z) = a \qquad z = \frac{x-\mu}{a}$$

$$p(x) = p(z)\left|\frac{1}{f'(z)}\right| = \frac{1}{\sqrt{2\pi}}e^{-\frac{1}{2}z^2}\frac{1}{a} = \frac{1}{\sqrt{2\pi a^2}}e^{-\frac{(x-\mu)^2}{2a^2}} = N(x|\mu, a^2)$$

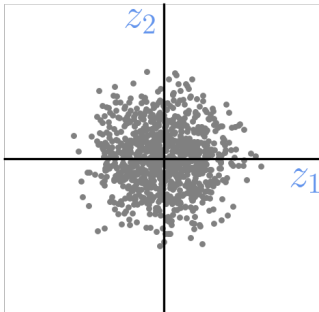**In general** $\quad \left.\begin{array}{l} z \sim N(\mu_z, \sigma_z^2) \\ x = az + b \end{array}\right] \quad \left.\begin{array}{l} E[x] = aE[z] + b = a\mu_z + b \\ \text{var}[x] = a^2\text{var}[z] = a^2\sigma_z^2 \end{array}\right] \quad x \sim N(a\mu_z + b, a^2\sigma_z^2)$
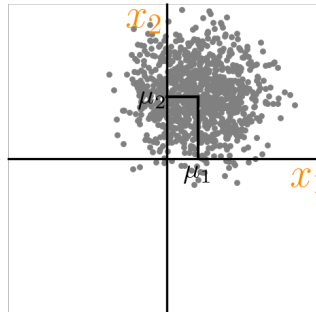
# Constructing multivariate Gaussians

**Changing mean**

**Changing mean & variance**

$$\left.\begin{array}{l} z_1 \sim \mathrm{N}(0,1) \\ z_2 \sim \mathrm{N}(0,1) \end{array}\right]\ z \sim \mathrm{N}(0, I)$$

$$\left.\begin{array}{l} x_1 = z_1 + \mu_1 \\ x_2 = z_2 + \mu_2 \end{array}\right]\ x = z + \mu$$

$$\left.\begin{array}{l} x_1 = a_{11} z_1 + a_{12} z_2 + \mu_1 \\ x_2 = a_{21} z_1 + a_{22} z_2 + \mu_2 \end{array}\right]\ x = Az + \mu$$



$$\mathrm{E}[z] = 0$$
$$\mathrm{cov}[z] = I$$

$$\mathrm{E}[x] = \mathrm{E}[z] + \mu = \mu$$
$$\mathrm{cov}[x] = \mathrm{cov}[z] = I$$

$$\mathrm{E}[x] = A\mathrm{E}[z] + \mu = \mu$$
$$\mathrm{cov}[x] = \text{see below} = AA^{\mathrm{T}}$$

$\mathrm{cov}[x]$ when changing mean & variance

$$x - E[x] = Az \qquad\qquad \mathrm{E}[zz^T] = \mathrm{cov}[z]$$

$$\mathrm{cov}[x] = \mathrm{E}[(x - \mathrm{E}[x])(x - \mathrm{E}[x])^T] = \mathrm{E}[Azz^T A^T] = A\mathrm{E}[zz^T]A^T = A\mathrm{cov}[z]A^T = AA^T$$

Also, **linear transformation of multivariate Gaussian remains Gaussian** (not shown)

**In general**

$$\left.\begin{array}{l} z \sim \mathrm{N}(\mu_z, \Sigma_z) \\ x = Az + \mu \end{array}\right]$$

$$\left.\begin{array}{l} \mathrm{E}[x] = A\mathrm{E}[z] + \mu = A\mu_z + \mu \\ \mathrm{cov}[x] = A\mathrm{cov}[z]A^T = A\Sigma_z A^T \end{array}\right]$$
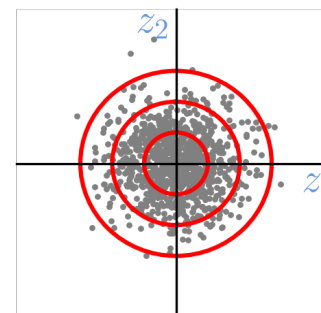
$$x \sim N(A\mu_z + \mu, A\Sigma_z A^T)$$

# Exploring covariance structure through isoprobability contours

**Standard Gaussian**

$$z \sim N(\mathbf{0}, \mathbf{I})$$

$$p(z) \propto e^{-\frac{1}{2}z^T z} = \text{const.} \longrightarrow z^T z = \text{const.}$$

$$\text{in 2D: } z_1^2 + z_2^2 = \text{const.}$$

**Introducing scaling $D$ (diagonal)**

$$y = Dz \qquad \text{cov}[y] = \Sigma_y = DD^T = D^2$$

$$y^T \Sigma_y^{-1} y = y^T D^{-2} y = z^T z = \text{const.}$$

$$\text{in 2D: } \left(\frac{y_1}{d_1}\right)^2 + \left(\frac{y_2}{d_2}\right)^2 = \text{const.}$$

$$D = \begin{bmatrix} 1.5 & 0 \\ 0 & 0.5 \end{bmatrix}$$

**Introducing rotation $R$ (orthonormal)**

$$x = Ry = RDz \qquad \text{cov}[x] = \Sigma_x = RDD^T R^T = RD^2 R^T$$

$$x^T \Sigma_x^{-1} x = x^T R D^{-2} R^T x = z^T z = \text{const.}$$

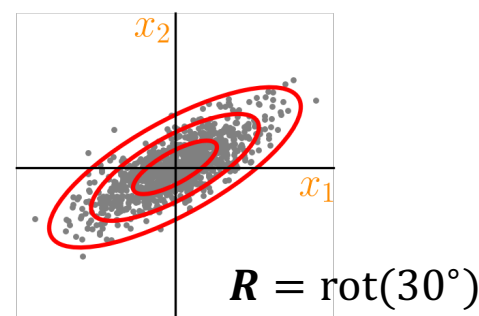$$\text{in 2D: } \left(\frac{(R^T x)_1}{d_1}\right)^2 + \left(\frac{(R^T x)_2}{d_2}\right)^2 = \text{const.}$$

$$R = \text{rot}(30°)$$

columns = $\Sigma_x$ eigenvectors $\qquad$ diagonal = (positive) $\Sigma_x$ eigenvalues

Spectral decomposition: $\Sigma_x = RD^2 R^T$

# Rank-deficient covariance matrices

So far: (implicitly) assumed $\dim(\boldsymbol{x}) \equiv N_x = N_z \equiv \dim(\boldsymbol{z})$ (i.e., a square $\boldsymbol{A}$)

$$N_x{\times}N_z \qquad N_z{\times}N_z$$
$$\boldsymbol{\Sigma}_x = \boldsymbol{R}\boldsymbol{D}^2\boldsymbol{R}^T$$

**Case 1**: Mapping from higher-dimensional $\boldsymbol{z}$ to lower-dimensional $\boldsymbol{x}$, $N_x < N_z$ (more "inputs" than "outputs")

$\boldsymbol{\Sigma}_x$ is full rank (i.e., all eigenvalues are non-zero)

**Case 2**: Mapping from lower-dimensional $\boldsymbol{z}$ to higher-dimensional $\boldsymbol{x}$, $N_x > N_z$ (fewer "inputs" than "outputs")

$\boldsymbol{\Sigma}_x$ is full rank-deficient (i.e., some eigenvalues are zero)

With $N_z = 1$ and $N_x = 2$: $\boldsymbol{x} = \boldsymbol{a}\boldsymbol{z} + \boldsymbol{\mu}$

# Interim summary: the Gaussian distribution

Fully specified by mean (vector) $\boldsymbol{\mu}$ and (co)variance (matrix) $\boldsymbol{\Sigma}$

Probability drops off quickly with distance from mean ("light" tail)
- sensitive to outliers

Linear transformation of (univariate/multivariate) Gaussian = Gaussian

Multivariate Gaussian = linear transformation of standard Gaussians
- Linear transformation = shift + scaling & rotation
- Scaling/rotation = eigenvalues/vectors of covariance matrix

Linear transformation from low-d to high-d: rank-deficient covariance (perfect correlations)

# Probabilistic operations: marginalization

Linear construction of multivariate Gaussian

$$z \sim N(0, I) \qquad x = Az + \mu \qquad x \sim N(\mu, \Sigma) \qquad \Sigma = AA^T$$

**Moments** of $x$ components

$i$th row of $A$

$$x_i = a_{i:}^T z + \mu_i$$

$$E[x_i] = \mu_i$$

$$\text{cov}[x_i, x_j] = E\big[\underbrace{a_{i:}^T z \ z^T a_{j:}}_{x_i - E[x_i]}\big] = a_{i:}^T E[zz^T]a_{j:} = a_{i:}^T a_{j:} = \Sigma_{ij}$$



**Marginalization:** 'removing' $x_i$ from $x$ by

$$p(x') = p(x_1, \ldots, x_{i-1}, x_{i+1}, \ldots, x_N) = \int p(x)dx_i$$

does not change moments of $x_{j \neq i}$
(as moments of $x_j$ only depend on $\mu_j$ and $a_{j:}$)

$$x \sim N(\mu, \Sigma) \qquad\qquad x' \sim N(\mu', \Sigma')$$



(Alternative approach: define $x' = Bx$, where $B$ 'picks' a subset of $x$ components)

# Probabilistic operations: conditioning

Linear construction of multivariate Gaussian

$$\mathbf{z} \sim \mathrm{N}(0, \mathbf{I}) \qquad \mathbf{x} = \mathbf{A}\mathbf{z} + \boldsymbol{\mu} \qquad \mathbf{x} \sim \mathrm{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \qquad \boldsymbol{\Sigma} = \mathbf{A}\mathbf{A}^T$$
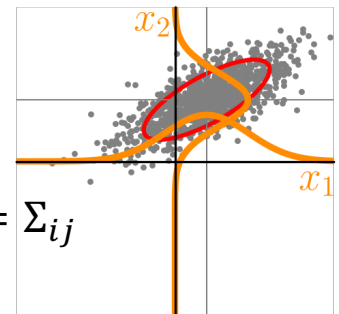
**Conditioning**: find $p(\mathbf{x}'|x_i) = p(x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_N | x_i)$ from $p(\mathbf{x})$

**To demonstrate**: assuming 2D $\mathbf{x}$

$x_1 = a_{11} z_1 + a_{12} z_2 + \mu_1$
$x_2 = a_{21} z_1 + a_{22} z_2 + \mu_2$



If we 'know' (i.e., condition on) $x_1$:

$$z_2 = \frac{x_1 - a_{11} z_1 - \mu_1}{a_{12}} \qquad \text{(the value of } z_2 \text{ is fully determined by } z_1 \text{ and } x_1)$$

$$x_2 = a_{21} z_1 + \frac{a_{22}(x_1 - a_{11} z_1 - \mu_1)}{a_{12}} + \mu_2$$

$$= \left(a_{21} - \frac{a_{22} a_{11}}{a_{12}}\right) z_1 + \frac{a_{22}(x_1 - \mu_1)}{a_{12}} + \mu_2 \qquad (x_1 \text{ is a linear function of } z_1 \rightarrow \text{Gaussian})$$

**More generally: conditionals of multivariate Gaussians are again Gaussian**

# Estimating parameters by maximum likelihood

Assume $N$ 1D observations $x_{1:N}$ drawn i.i.d. from Gaussian $N(x_n | \mu, \sigma^2)$

We find **maximum likelihood** parameters by solving

$$\hat{\mu}_{ML}, \hat{\sigma}^2_{ML} = \text{argmax}_{\mu, \sigma^2} \prod_{n=1}^{N} N(x_n | \mu, \sigma^2) = \text{argmax}_{\mu, \sigma^2} \sum_{n=1}^{N} \log N(x_n | \mu, \sigma^2)$$

$$= \text{argmax}_{\mu, \sigma^2} -\frac{N}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{n=1}^{N} (x_n - \mu)^2$$

Setting the derivative with respect to $\mu$ and $\sigma^2$ to zero yields

$$\hat{\mu}_{ML} = \frac{1}{N} \sum_{n} x_n \qquad \hat{\sigma}^2_{ML} = \frac{1}{N} \sum_{n} (x_n - \hat{\mu}_{ML})^2$$

(note: ML variance estimate is biased; see PRML Sec. 1.2.4)

**ML estimates are simply sample mean & sample variance**
(this also holds for multivariate Gaussian)

# Priors as implicitly added observations

**Gaussian prior** $p(\mu) = N\left(\mu \mid 0, \frac{\sigma^2}{\lambda}\right)$ **on mean** $\mu$ results in MAP estimate

$$\hat{\mu}_{MAP} = \text{argmax}_\mu\, p(\mu) \prod_{n=1}^{N} p(x_n \mid \mu, \sigma^2) = \frac{1}{N+\lambda}\left(\sum_{n=1}^{N} x_n + \lambda 0\right)$$

Prior implicitly adds $\lambda$ observations with value zero;
$\lambda \to 0$ recovers ML estimate, $\hat{\mu}_{MAP} \to \hat{\mu}_{ML}$; reveals implicit assumption of ML estimate

(Conjugate) **Normal Inverse Gamma prior** $p(\mu, \sigma^2) = N\left(\mu \mid 0, \frac{\sigma^2}{\lambda}\right) IG(\sigma^2 \mid \alpha, \beta)$ results in

$$\hat{\mu}_{MAP} = \frac{1}{N+\lambda} \sum_{n=1}^{N} x_n$$

Prior implicitly adds

$\lambda$ observations with value zero

$$\hat{\sigma}^2_{MAP} = \frac{2\alpha+3}{N+2\alpha+3}\frac{\lambda\hat{\mu}_{ML}+2\beta}{2\alpha+3} + \frac{N}{N+2\alpha+3}\hat{\sigma}^2_{ML}$$

$2\alpha + 3$ observations
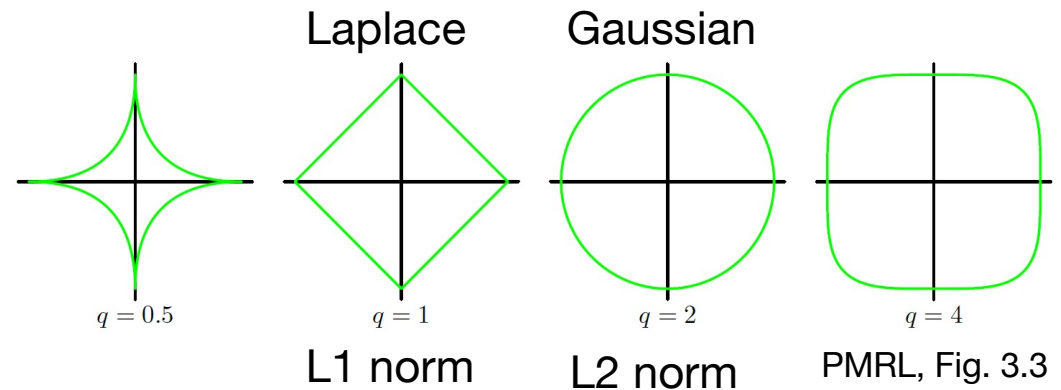with value $\frac{\lambda\hat{\mu}_{ML}+2\beta}{2\alpha+3}$

(Note: for the above, one can analytically find the whole posterior $p(\mu, \sigma^2 \mid x_{1:N})$)

# Priors as regularizers

MAP estimate balances prior $p(\boldsymbol{\mu})$ and likelihood $p(\boldsymbol{x}_{1:N}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$

$$\widehat{\boldsymbol{\mu}}_{MAP} = \text{argmax}_{\mu}\left(\log p(\boldsymbol{\mu}) + \log p(\boldsymbol{x}_{1:N}|\boldsymbol{\mu}, \boldsymbol{\Sigma})\right)$$
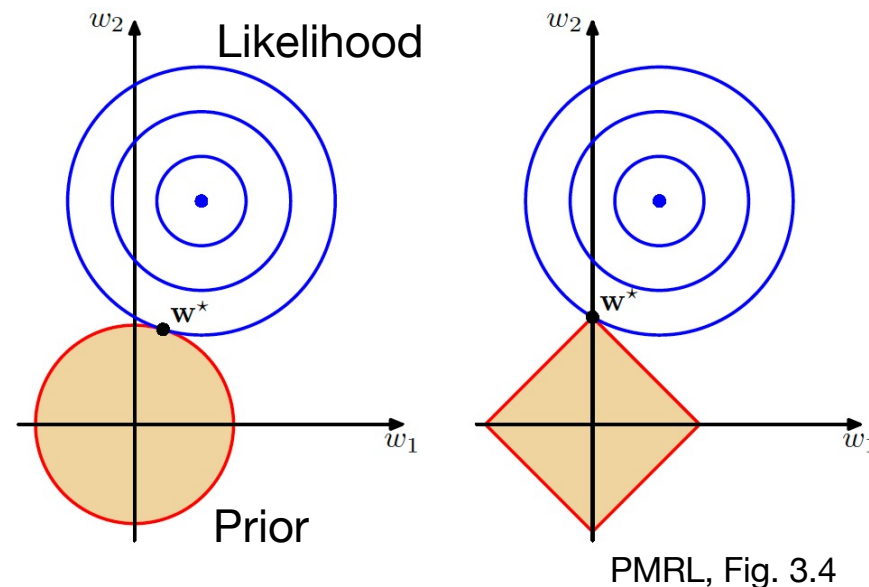
Assume parametric prior,
$\log p(\boldsymbol{\mu}) = \sum_{i=1}^{N_x} |\mu_i|^q + \text{const.}$

Laplace   Gaussian



$q = 0.5$     $q = 1$     $q = 2$     $q = 4$

L1 norm     L2 norm     PMRL, Fig. 3.3

Laplace prior pushes
MAP estimate
components to zero

*LASSO*:
Laplace prior = "LASSO"

*Elastic net*:
Gaussian + Laplace prior



Likelihood

Prior

PMRL, Fig. 3.4

# Interim summary: inference & priors

Both marginalization & conditioning of multivariate Gaussian yields Gaussian

Marginalization: removing rows / columns of $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$

Summary of inference on Gaussians: PRML Appendix B, Sec. *Gaussian*

Maximum likelihood = matching sample moments (biased)

Priors can be interpreted as additional observations contributing to posterior / MAP estimate

Priors regularize / Laplace prior (L1 regularization) pushes (some) MAP components to zero

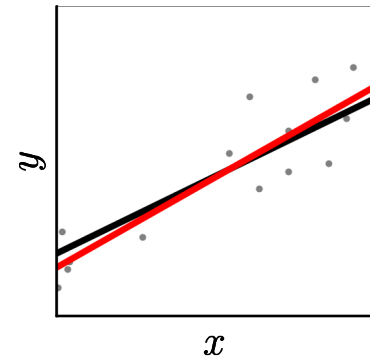# Overview

# Linear models: maximum likelihood = least squares

Assume 'output' $y$ is a linear function of 'input' $x$ + Gaussian noise

$$y = \boldsymbol{w}^T \boldsymbol{x} + w_0 + \eta \quad \leftarrow \eta \sim N(0, \sigma^2)$$

$$p(y|\boldsymbol{x}, \boldsymbol{w}) = N(y|\boldsymbol{w}^T \boldsymbol{x} + w_0, \sigma^2)$$



**Linear regression**

Observe $\{\boldsymbol{x}_n, y_n\}_{n=1}^N$, what is $\boldsymbol{w}$ and $w_0$ that makes data most likely?

$$\text{argmax}_{\boldsymbol{w}, w_0} \log \prod_{n=1}^N p(y_n|\boldsymbol{x}_n, \boldsymbol{w}) = \text{argmax}_{\boldsymbol{w}, w_0} - \frac{1}{2\sigma^2} \sum_{n=1}^N (\underbrace{y_n}_{\text{target}} - \underbrace{\boldsymbol{w}^T \boldsymbol{x}_n - w_0}_{\text{model}})^2$$

**Maximizing data likelihood = minimizing squared error**
(sensitive to outliers)

**Solution** (assuming $w_0 = 0$, we'll get to this)

Setting gradient of LLH with respect to $\boldsymbol{w}$ to zero

$$\widehat{\boldsymbol{w}}_{ML} = \left(\sum_n \boldsymbol{x}_n \boldsymbol{x}_n^T\right)^{-1} \sum_n \boldsymbol{x}_n y_n = (\boldsymbol{X}^T \boldsymbol{X})^{-1} \boldsymbol{X}^T \boldsymbol{y}$$

*design matrix*

$$\boldsymbol{X} = \begin{pmatrix} - \boldsymbol{x}_1^T - \\ - \boldsymbol{x}_2^T - \\ \vdots \end{pmatrix} \quad \boldsymbol{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \end{pmatrix}$$

# Linear models: linear in parameters, not 'inputs'

'Inputs' are assumed known/given – can change them while keeping solution structure

$$y = \boldsymbol{w}^T \boldsymbol{\phi}(\boldsymbol{x}) + \eta \leftarrow \eta \sim \mathrm{N}(0, \sigma^2)$$

$$p(y|\boldsymbol{x}, \boldsymbol{w}) = \mathrm{N}(y|\boldsymbol{w}^T \boldsymbol{\phi}(\boldsymbol{x}), \sigma^2)$$

Solution keep structure: $\widehat{\boldsymbol{w}}_{ML} = (\boldsymbol{\Phi}^T \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^T \boldsymbol{y}$

$$\boldsymbol{\Phi} = \begin{pmatrix} -\boldsymbol{\phi}(\boldsymbol{x}_1)^T - \\ -\boldsymbol{\phi}(\boldsymbol{x}_2)^T - \\ \vdots \end{pmatrix}$$

**Example:** include bias term (assumed from now on)

$$\boldsymbol{\phi}(\boldsymbol{x}) = \begin{pmatrix} 1 \\ \boldsymbol{x} \end{pmatrix} \qquad \boldsymbol{w} = \begin{pmatrix} w_0 \\ \widetilde{\boldsymbol{w}} \end{pmatrix} \longrightarrow \boldsymbol{w}^T \boldsymbol{\phi}(\boldsymbol{x}) = \widetilde{\boldsymbol{w}}^T \boldsymbol{x} + w_0$$

**Example:** polynomial regression

$$\boldsymbol{\phi}(x) = \begin{pmatrix} 1 \\ x \\ x^2 \\ \vdots \end{pmatrix} \longrightarrow \boldsymbol{w}^T \boldsymbol{\phi}(x) = w_0 + w_1 x + w_2 x^2 + \cdots$$

# Linear models: revisiting (parameter) priors

$$\widehat{w}_{MAP} = \text{argmax}_w\, p(\boldsymbol{y}|\boldsymbol{w}, \boldsymbol{X})p(\boldsymbol{w}) = \text{argmax}_w\, (\log p(\boldsymbol{y}|\boldsymbol{w}, \boldsymbol{X}) + \log p(\boldsymbol{w}))$$

Assume parametric prior,
$$\log p(\boldsymbol{w}) = \sum_{i=1}^{N_x} |w_i|^q + \text{const.}$$

Laplace    Gaussian



$q = 0.5$     $q = 1$     $q = 2$     $q = 4$

L1 norm    L2 norm    PMRL, Fig. 3.3

Laplace prior pushes MAP estimate
components to zero – some 'inputs'
get ignored / don't matter

*LASSO*:
Laplace prior = "LASSO"

*Elastic net*:
Gaussian + Laplace prior



PMRL, Fig. 3.4

# Bayesian linear regression

posterior ∝ likelihood x prior

$$p(\boldsymbol{w}|\boldsymbol{X}, \boldsymbol{y}) \propto p(\boldsymbol{y}|\boldsymbol{w}, \boldsymbol{X})p(\boldsymbol{w})$$

Gaussian in $\boldsymbol{w}$



PRML Fig 3.7

# Interim summary: linear models

Linear models are linear in the parameters

Using non-linear functions of 'inputs' increases linear model applicability

Simple analytical solutions for ML estimates,
and MAP estimates / full posterior with Gaussian priors

ML solution = least squares solution (due to Gaussian residuals/noise)

L1 / Laplace parameter prior introduces sparsity

# Overview

# Bayesian model comparison

So far implicit: models as yet another variable to condition on

e.g., linear function     e.g., quadratic function

Assume models $M_1$ and $M_2$ with priors $p(\mathbf{w}|M_j)$ and data likelihood $p(\mathbf{X}|\mathbf{w}, M_j)$

Observe $\mathbf{X}$: probability of model $M_j$ is

$$p(M_j|\mathbf{X}) \propto p(\mathbf{X}|M_j)p(M_j)$$

**Model evidence** $p(\mathbf{X}|M_j) = \int p(\mathbf{X}|\mathbf{w}, M_j)p(\mathbf{w}|M_j)dw$

**Bayes Factor:** comparing models

$p(M_1|\mathbf{X}) > p(M_2|\mathbf{X})$, and $p(M_1) = p(M_2)$, implies

Bayes Factor $K \equiv \dfrac{p(\mathbf{X}|M_1)}{p(\mathbf{X}|M_2)} = \dfrac{p(M_1|\mathbf{X})}{p(M_2|\mathbf{X})} > 1$

| $K$ | dHart | bits | Strength of evidence |
|---|---|---|---|
| $< 10^0$ | $< 0$ | $< 0$ | Negative (supports $M_2$) |
| $10^0$ to $10^{1/2}$ | 0 to 5 | 0 to 1.6 | Barely worth mentioning |
| $10^{1/2}$ to $10^1$ | 5 to 10 | 1.6 to 3.3 | Substantial |
| $10^1$ to $10^{3/2}$ | 10 to 15 | 3.3 to 5.0 | Strong |
| $10^{3/2}$ to $10^2$ | 15 to 20 | 5.0 to 6.6 | Very strong |
| $> 10^2$ | $> 20$ | $> 6.6$ | Decisive |

Jeffreys (1998); Wikipedia

# Balancing model fit and complexity

Bayes Factor relies on $p(X|M_j)$: how likely is it to observe data $X$ under given model $M_j$?

$\int p(X|M_j)dX = 1$ $\longrightarrow$ the more different $X$'s a model can 'explain' (i.e., $p(X|M_j) \gg 0$), the lower $p(X|M_j)$ for each $X$

**Example**: linear vs. quadratic function



all quadratic functions
$y = w_0 + w_1 x + w_2 x^2 + \eta$

all linear function
$y = w_0 + w_1 x + \eta$

$+ X_{quad}$

$+ X_{lin}$

$p(X_{lin}|M_{lin}) >$
$p(X_{lin}|M_{quad})$

$p(X_{quad}|M_{lin}) \approx 0$

# Adjusting model complexity through hyperpriors

Hyperpriors: priors $p(\boldsymbol{\alpha})$ on parameters $\boldsymbol{\alpha}$ of the prior $p(\boldsymbol{w}|\boldsymbol{\alpha})$

**Example**: hyperprior on quadratic coefficient

$$p(w_2|\,\alpha) = N(w_2|0, \alpha^{-1})$$

$\alpha \to 0 \quad w_2$ can take arbitrary values
$\alpha \to \infty \quad w_2 \to 0$ (linear function)

$$y = w_0 + w_1 x + w_2 x^2 + \eta$$

Inferring $\alpha$ from data: are quadratic or linear models more adequate?

$\longrightarrow$ determines model complexity from data
(a form of Automated Relevance Determination; ARD)

**Full inference** (usually intractable):

posterior predictive density          data likelihood

$$p(\boldsymbol{x}|\boldsymbol{X}) = \iint p(\boldsymbol{x}|\boldsymbol{w})p(\boldsymbol{w}|\boldsymbol{\alpha}, \boldsymbol{X})p(\boldsymbol{\alpha}|\boldsymbol{X})\mathrm{d}\boldsymbol{w}\mathrm{d}\boldsymbol{\alpha}$$

posterior          hyperposterior

$$p(\boldsymbol{\alpha}|\boldsymbol{X}) \propto p(\boldsymbol{X}|\boldsymbol{\alpha})p(\boldsymbol{\alpha})$$

$$= \int p(\boldsymbol{X}|\boldsymbol{w})p(\boldsymbol{w}|\boldsymbol{\alpha})\mathrm{d}\boldsymbol{w}$$

# Evidence approximation

Tractable hyperparameter inference

$$p(\boldsymbol{x}|\boldsymbol{X}) = \iint p(\boldsymbol{x}|\boldsymbol{w})p(\boldsymbol{w}|\boldsymbol{\alpha},\boldsymbol{X})p(\boldsymbol{\alpha}|\boldsymbol{X})\mathrm{d}\boldsymbol{w}\mathrm{d}\boldsymbol{\alpha}$$

hyperposterior

<span style="color:red">approximate by MAP, $\widehat{\boldsymbol{\alpha}}_{MAP}$</span>

$$p(\boldsymbol{x}|\boldsymbol{X}) \approx p(\boldsymbol{x}|\boldsymbol{X},\widehat{\boldsymbol{\alpha}}_{MAP}) = \int p(\boldsymbol{x}|\boldsymbol{w})p(\boldsymbol{w}|\widehat{\boldsymbol{\alpha}}_{MAP},\boldsymbol{X})\mathrm{d}\boldsymbol{w}$$

**Estimating MAP hyperparameters**

$$\widehat{\boldsymbol{\alpha}}_{MAP} = \mathrm{argmax}_{\boldsymbol{\alpha}}\log p(\boldsymbol{\alpha}|\boldsymbol{X}) = \mathrm{argmax}_{\boldsymbol{\alpha}}(\log p(\boldsymbol{X}|\boldsymbol{\alpha}) + \log p(\boldsymbol{\alpha}))$$

$$p(\boldsymbol{X}|\boldsymbol{\alpha}) = \int p(\boldsymbol{X}|\boldsymbol{w})p(\boldsymbol{w}|\boldsymbol{\alpha})\mathrm{d}\boldsymbol{w}$$

In PRML, Ch. 3.5:
- Assume uninformative $p(\boldsymbol{\alpha})$: $\widehat{\boldsymbol{\alpha}}_{MAP} \approx \widehat{\boldsymbol{\alpha}}_{ML}$
- $\widehat{\boldsymbol{\alpha}}_{MAP}$ depends on mean estimate, which depends on $\widehat{\boldsymbol{\alpha}}_{MAP}$: estimate recursively

# Interim summary: Bayesian model comparison

Comparing models by marginalizing over the parameters / hyperparameters

Implicit complexity penalty by lower model evidence for models that can fit more data

Hyperpriors can provide automatic inference of desired model complexity

Evidence approximation for tractable hyperparameter marginalization

# Overview

**Univariate and multivariate Gaussian distributions**

Means and variances of linearly transformed random variables (general)

Univariate standard Gaussian and linear transformations thereof

Multivariate Gaussians through linear constructions

Probabilistic operations: marginalization and conditioning

Maximum likelihood parameter estimation

Priors as added observations / regularizers

**Linear models**

Maximum likelihood and least squares estimates

What is linear in linear models?

Priors and Bayesian inference

**Bayesian model comparison**

Trading off goodness-of-fit with model complexity

Adjusting model complexity through hyperpriors

Evidence approximation

# Summary

Gaussians are fully described by their mean (vector) and (co)variance matrix

A linear transformation of a Gaussian remains a Gaussian

Marginalization and conditioning Gaussians (on Gaussians) yields Gaussians

Priors can act as additional observations / as regularizers


Linear models: linear in parameters, not 'inputs'

Maximum likelihood in linear models = least squares, with analytic solution


Bayesian model comparison trades off goodness-of-fit with model complexity

Hyperparameters can continuously modulate model complexity (evidence approx.)

# Until next week

Read paper and prepare presentation (see notes for Session 3)

Read statistical methods sections (see notes for Session 3)

**Next session**

Q&A for previous session (~15min)

Paper discussions (~1h)

Introducing generalized linear models (~30min)