

CS189: Intro to Machine Learning

Summer 2018

Lecture 5: Bias-variance tradeoff and multivariate gaussians

Josh Tobin
UC Berkeley EECS

Announcements

- HW2 posted

Outline for today

- Bias-variance decomposition
- Multivariate gaussians

Outline for today

- **Bias-variance decomposition**
- Multivariate gaussians

Bias-variance decomposition

$$\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$$

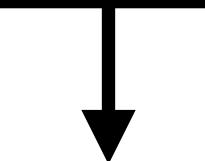
View x_i as sampled from r.v. X

Good model: one with low *expected* error

Bias-variance tradeoff

Goal: We want to find the model family with the lowest *expected* error

$$\epsilon(\vec{x}; f) = \mathbb{E} [(f_\theta^*(\vec{x}; \mathcal{D}) - y)^2] = \mathbb{E}_{y, \mathcal{D}} [(f_\theta^*(\vec{x}; \mathcal{D}) - y)^2]$$

$$(f_\theta^*(\vec{x}; \mathcal{D}) - y)^2$$


“Best” model: choice of theta
with the lowest error on D

Bias-variance tradeoff

Goal: We want to find the model family with the lowest *expected* error

$$\epsilon(\vec{x}; f) = \mathbb{E} [(f_\theta^*(\vec{x}; \mathcal{D}) - y)^2] = \mathbb{E}_{y, \mathcal{D}} [(f_\theta^*(\vec{x}; \mathcal{D}) - y)^2]$$

$$\frac{(f_\theta^*(\vec{x}; \mathcal{D}) - y)^2}{\downarrow}$$

“Best” model’s estimate for y at
the datapoint x

Bias-variance tradeoff

Goal: We want to find the model family with the lowest *expected* error

$$\epsilon(\vec{x}; f) = \mathbb{E} [(f_\theta^*(\vec{x}; \mathcal{D}) - y)^2] = \mathbb{E}_{y, \mathcal{D}} [(f_\theta^*(\vec{x}; \mathcal{D}) - y)^2]$$

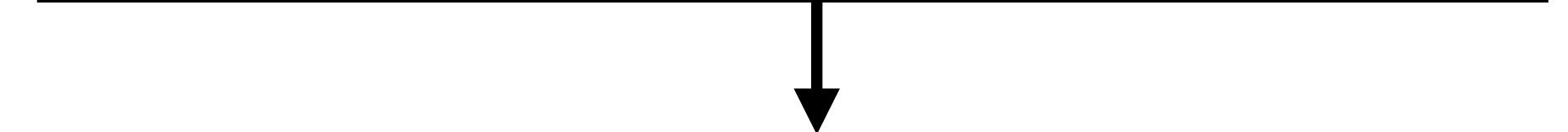
$$\frac{(f_\theta^*(\vec{x}; \mathcal{D}) - y)^2}{\downarrow}$$

Squared error between estimate for y and one choice for the observed y at datapoint x

Bias-variance tradeoff

Goal: We want to find the model family with the lowest *expected* error

$$\epsilon(\vec{x}; f) = \mathbb{E} [(f_\theta^*(\vec{x}; \mathcal{D}) - y)^2] = \mathbb{E}_{y, \mathcal{D}} [(f_\theta^*(\vec{x}; \mathcal{D}) - y)^2]$$

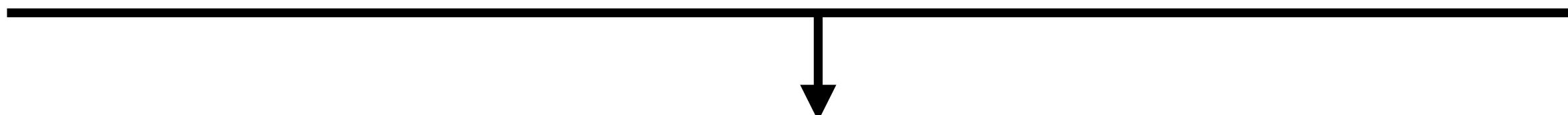
$$\frac{\mathbb{E}_{y \sim p(y|\vec{x})} [(f_\theta^*(\vec{x}; \mathcal{D}) - y)^2]}{p(y \mid \vec{x}) = \mathcal{N}(f(x), \sigma^2)}$$


Average squared error between estimate
for y and true y

Bias-variance tradeoff

Goal: We want to find the model family with the lowest *expected* error

$$\epsilon(\vec{x}; f) = \mathbb{E} [(f_\theta^*(\vec{x}; \mathcal{D}) - y)^2] = \mathbb{E}_{y, \mathcal{D}} [(f_\theta^*(\vec{x}; \mathcal{D}) - y)^2]$$

$$\frac{\mathbb{E}_{\mathcal{D} \sim p(\mathcal{D})} \mathbb{E}_{y \sim p(y|\vec{x})} [(f_\theta^*(\vec{x}; \mathcal{D}) - y)^2]}{p(\mathcal{D})} \quad p(\mathcal{D}) = \prod_{i=1}^n p(x_i)p(y_i \mid x_i)$$


Average over all datasets of mean squared error on that dataset

Bias-variance tradeoff

Goal: We want to find the model family with the lowest *expected* error

$$\epsilon(\vec{x}; f) = \mathbb{E} [(f_\theta^*(\vec{x}; \mathcal{D}) - y)^2] = \mathbb{E}_{y, \mathcal{D}} [(f_\theta^*(\vec{x}; \mathcal{D}) - y)^2]$$

$$\mathbb{E}_{\mathcal{D} \sim p(\mathcal{D})} \mathbb{E}_{y \sim p(y|\vec{x})} [(f_\theta^*(\vec{x}; \mathcal{D}) - y)^2] = \mathbb{E} [(f_\theta^*(\vec{x}; \mathcal{D}) - y)^2]$$

Bias-variance tradeoff

Goal: We want to find the model family with the lowest *expected* error

$$\epsilon(\vec{x}; f) = \mathbb{E}_{y, \mathcal{D}} [(f_\theta^*(\vec{x}; \mathcal{D}) - y)^2]$$

Let's decompose this into parts

Bias-variance tradeoff

Goal: We want to find the model family with the lowest *expected* error

$$\epsilon(\vec{x}; f) = \mathbb{E}_{y, \mathcal{D}} [(f_\theta^*(\vec{x}; \mathcal{D}) - y)^2]$$

Let's decompose this into parts

$$= (\mathbb{E}[f_\theta^*(\vec{x}; \mathcal{D})] - \mathbb{E}[y])^2 + \text{Var}(f_\theta^*(\vec{x}; \mathcal{D})) + \text{Var}(z)$$

The equation is decomposed into three parts by horizontal lines and arrows pointing downwards. The first part, $(\mathbb{E}[f_\theta^*(\vec{x}; \mathcal{D})] - \mathbb{E}[y])^2$, is labeled "bias^2". The second part, $\text{Var}(f_\theta^*(\vec{x}; \mathcal{D}))$, is labeled "Variance". The third part, $\text{Var}(z)$, is labeled "Irreducible error".

bias²

Variance Irreducible
error

Bias-variance tradeoff

$$= (\mathbb{E} [f_\theta^*(\vec{x}; \mathcal{D})] - \mathbb{E} [y])^2 + \text{Var}(f_\theta^*(\vec{x}; \mathcal{D})) + \text{Var}(z)$$

bias²

Variance

Irreducible
error

$$(\mathbb{E} [f_\theta^*(\vec{x}; \mathcal{D})] - \mathbb{E} [y])^2$$

Average over all datasets of that dataset's
best model's estimate of $f(x)$

Bias-variance tradeoff

$$= (\mathbb{E}[f_\theta^*(\vec{x}; \mathcal{D})] - \mathbb{E}[y])^2 + \text{Var}(f_\theta^*(\vec{x}; \mathcal{D})) + \text{Var}(z)$$

bias²

Variance

Irreducible
error

$$(\mathbb{E}[f_\theta^*(\vec{x}; \mathcal{D})] - \mathbb{E}[y])^2$$

Average of noisy y's for a given x

Bias-variance tradeoff

$$= (\mathbb{E}[f_\theta^*(\vec{x}; \mathcal{D})] - \mathbb{E}[y])^2 + \text{Var}(f_\theta^*(\vec{x}; \mathcal{D})) + \text{Var}(z)$$

bias²

Variance

Irreducible
error

$$(\mathbb{E}[f_\theta^*(\vec{x}; \mathcal{D})] - \mathbb{E}[y])^2$$

Average error in estimate of y over all
datasets

Bias-variance tradeoff

$$= (\mathbb{E}[f_\theta^*(\vec{x}; \mathcal{D})] - \mathbb{E}[y])^2 + \text{Var}(f_\theta^*(\vec{x}; \mathcal{D})) + \text{Var}(z)$$

bias²

Variance

Irreducible
error

$$\text{Var}(f_\theta^*(\vec{x}; \mathcal{D}))$$

Variance of estimate of y over all possible datasets

Bias-variance tradeoff

$$= (\mathbb{E}[f_\theta^*(\vec{x}; \mathcal{D})] - \mathbb{E}[y])^2 + \text{Var}(f_\theta^*(\vec{x}; \mathcal{D})) + \text{Var}(z)$$

bias²

Variance

Irreducible
error

Var(z)

Variance of the noise term
in the model for y

Derive bias-variance decomposition

Goal: We want to find the model family with the lowest *expected* error

$$\epsilon(\vec{x}; f) = \mathbb{E} [(f_{\theta}^*(\vec{x}; \mathcal{D}) - y)^2]$$

Start by analyzing mean & var of y

$$\mathbb{E}[Y] = \mathbb{E}[f(\vec{x}) + z] = f(\vec{x}) + \mathbb{E}[z] = f(\vec{x})$$

$$\text{Var}(Y) = \text{Var}(f(\vec{x}) + z) = \text{Var}(Z)$$

Bias-variance tradeoff

Goal: We want to find the model family with the lowest *expected* error

$$\epsilon(\vec{x}; f) = \mathbb{E} [(f_\theta^*(\vec{x}; \mathcal{D}) - y)^2]$$

Add one more lemma (for any RV)

$$\text{Var}(X) = \mathbb{E} [(X - \mathbb{E}(X))^2] = \mathbb{E} [X^2] - \mathbb{E} [X]^2$$

$$\implies \mathbb{E} [X^2] = \text{Var}(X) + \mathbb{E} [X]^2$$

Bias-variance tradeoff

$$\begin{aligned}\epsilon(\vec{x}; f) &= \mathbb{E} [(f_\theta^*(\vec{x}; \mathcal{D}) - y)^2] \\&= \mathbb{E}(f_\theta^*(\vec{x}; \mathcal{D})^2) + \mathbb{E}(y^2) - 2\mathbb{E}[f_\theta^*(\vec{x}; \mathcal{D}) \cdot y] \\&= (\text{Var}(f_\theta^*(\vec{x}; \mathcal{D})) + \mathbb{E}[f_\theta^*(\vec{x}; \mathcal{D})]^2) \\&\quad + (\text{Var}(Y) + \mathbb{E}[y]^2) - 2\mathbb{E}[f_\theta^*(\vec{x}; \mathcal{D})] \cdot \mathbb{E}[y] \\&= (\mathbb{E}[f_\theta^*(\vec{x}; \mathcal{D})]^2 - 2\mathbb{E}[f_\theta^*(\vec{x}; \mathcal{D})] \cdot \mathbb{E}[y] + \mathbb{E}[y]^2) \\&\quad + \text{Var}(f_\theta^*(\vec{x}; \mathcal{D})) + \text{Var}(z) \\&= (\mathbb{E}[f_\theta^*(\vec{x}; \mathcal{D})] - \mathbb{E}[y])^2 + \text{Var}(f_\theta^*(\vec{x}; \mathcal{D})) + \text{Var}(z)\end{aligned}$$

Facts

$$\mathbb{E}[X^2] = \text{Var}(X) + \mathbb{E}[X]^2$$

$$\text{Var}(Y) = \text{Var}(z)$$

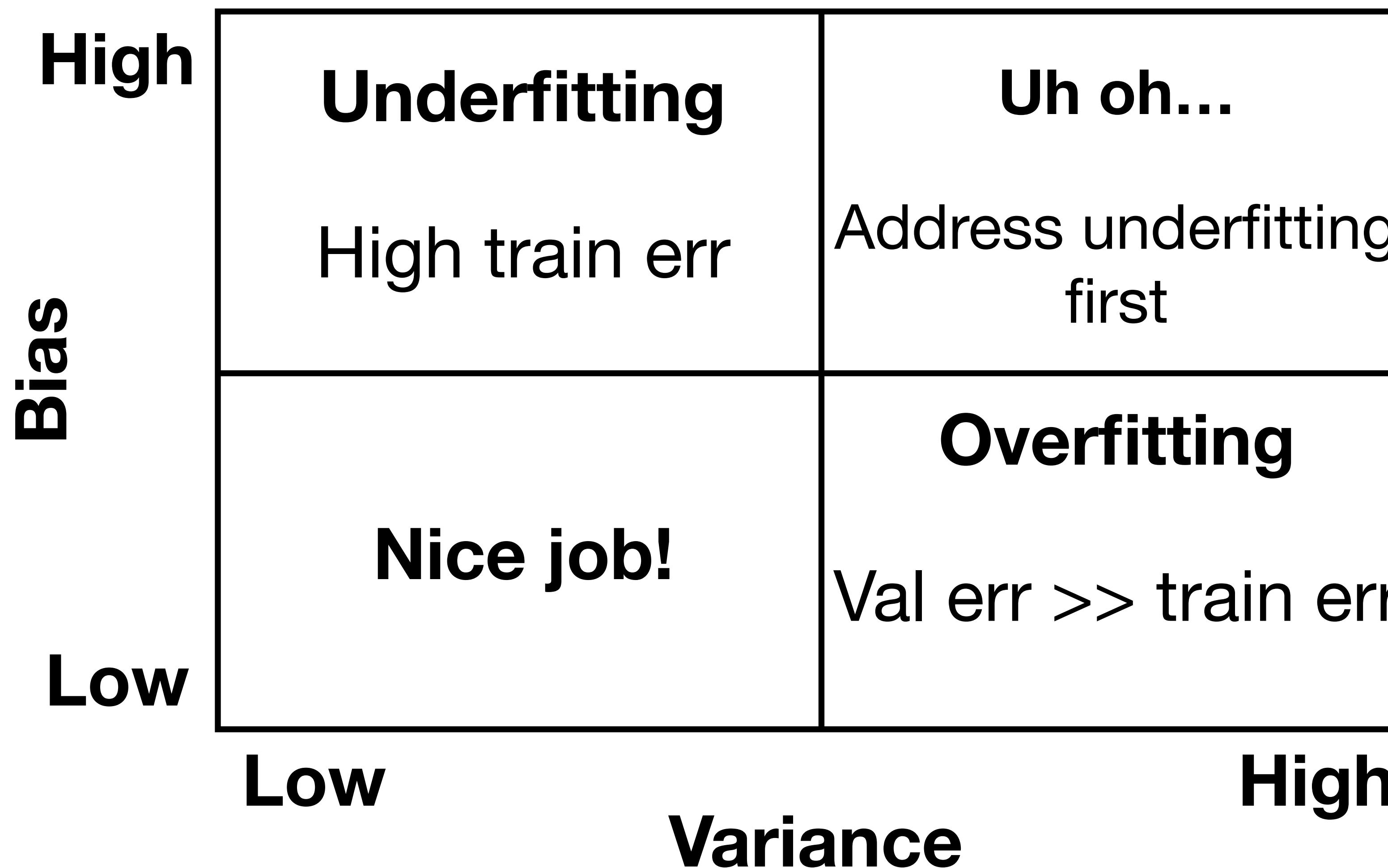
$$\mathbb{E}(Y) = f(\vec{x})$$

Bias-variance tradeoff

$$\epsilon(\vec{x}; f) = \mathbb{E}_{\vec{x}, y, \mathcal{D}} [(f_\theta^*(\vec{x}; \mathcal{D}) - y)^2]$$
$$= \underbrace{(\mathbb{E} [f_\theta^*(\vec{x}; \mathcal{D})] - \mathbb{E} [y])^2}_{\text{bias}^2} + \underbrace{\text{Var}(f_\theta^*(\vec{x}; \mathcal{D}))}_{\text{Variance}} + \underbrace{\text{Var}(z)}_{\text{Irreducible error}}$$

Bias-variance cheat sheet

$$\epsilon(\vec{x}; f) = (\mathbb{E} [f_\theta^*(\vec{x}; \mathcal{D})] - \mathbb{E} [y])^2 + \text{Var}(f_\theta^*(\vec{x}; \mathcal{D})) + \text{Var}(z)$$



How to trade off between bias and variance?

| Category | Intervention | Bias | Variance |
|----------------|---|------|----------|
| Data | <ul style="list-style-type: none">• Add data | | |
| Features | <ul style="list-style-type: none">• Add good features | | |
| | <ul style="list-style-type: none">• Add bad features | | |
| | <ul style="list-style-type: none">• Remove features | | |
| Regularization | <ul style="list-style-type: none">• Increase regularization | | |
| | <ul style="list-style-type: none">• Decrease regularization | | |

Outline for today

- Bias-variance decomposition
- **Review of multivariate gaussians**

Breaking down the iid assumptions

Recall:

$$y_i = f(\vec{x}_i) + z_i \quad z_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2)$$

Where does iid fail?

- Time series
- “Sliding windows”
- Etc

Multivariate gaussian distribution

Definition

A joint distribution over $\{x_1, \dots, x_k\}$
with the following pdf:

$$p(x_1, \dots, x_k) =$$

Multivariate gaussian distribution

Definition

A joint distribution over $\{x_1, \dots, x_k\}$
with the following pdf:

$$p(x_1, \dots, x_k) = \frac{\exp(\cdot)}{(\cdot)}$$

Multivariate gaussian distribution

Definition

A joint distribution over $\{x_1, \dots, x_k\}$
with the following pdf:

$$p(x_1, \dots, x_k) = \frac{\exp\left(-\frac{1}{2}(\vec{x} - \vec{\mu})^T \Sigma^{-1} (\vec{x} - \vec{\mu})\right)}{(\cdot)}$$

Multivariate gaussian distribution

Definition

A joint distribution over $\{x_1, \dots, x_k\}$
with the following pdf:

$$p(x_1, \dots, x_k) = \frac{\exp\left(-\frac{1}{2}(\vec{x} - \vec{\mu})^T \Sigma^{-1} (\vec{x} - \vec{\mu})\right)}{\sqrt{(2\pi)^k |\Sigma|}}$$

Equivalent definition

Definition

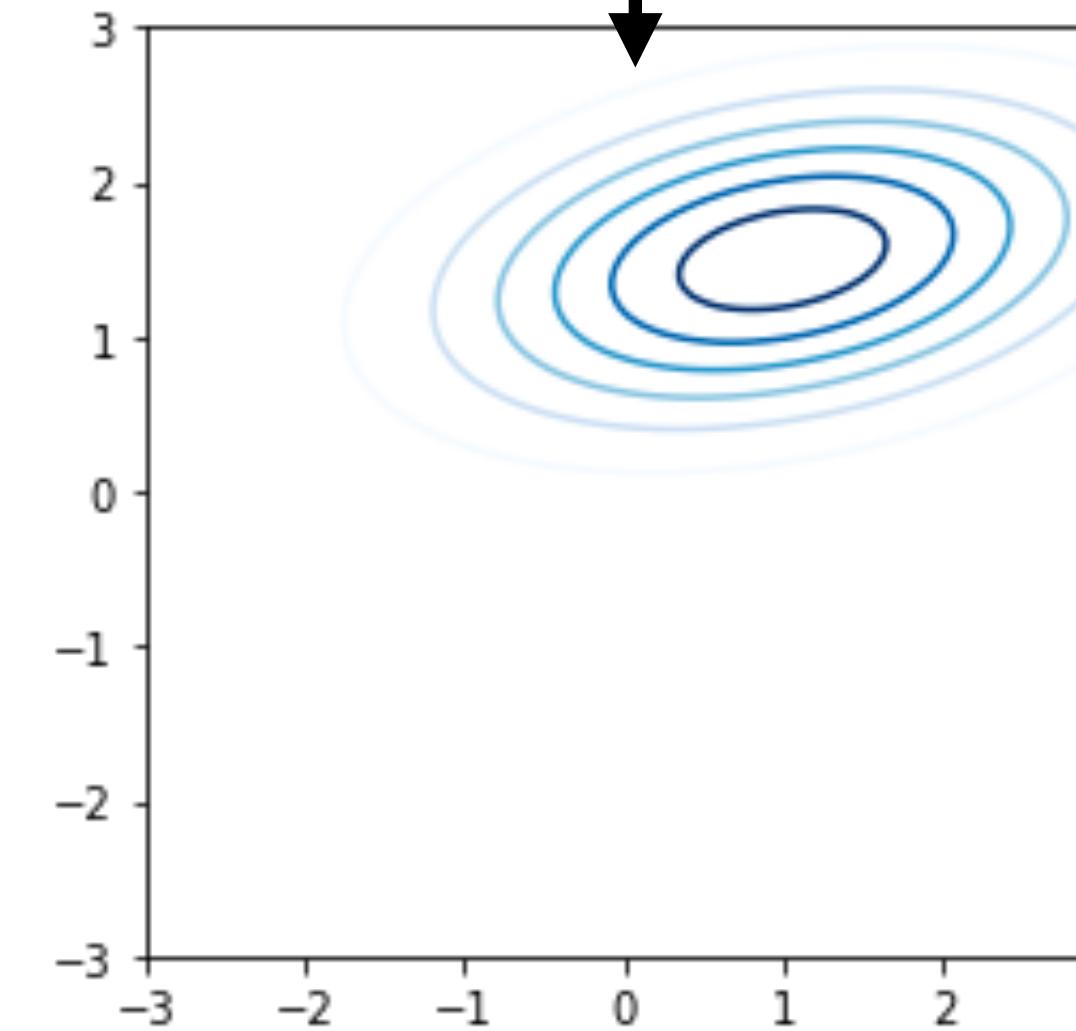
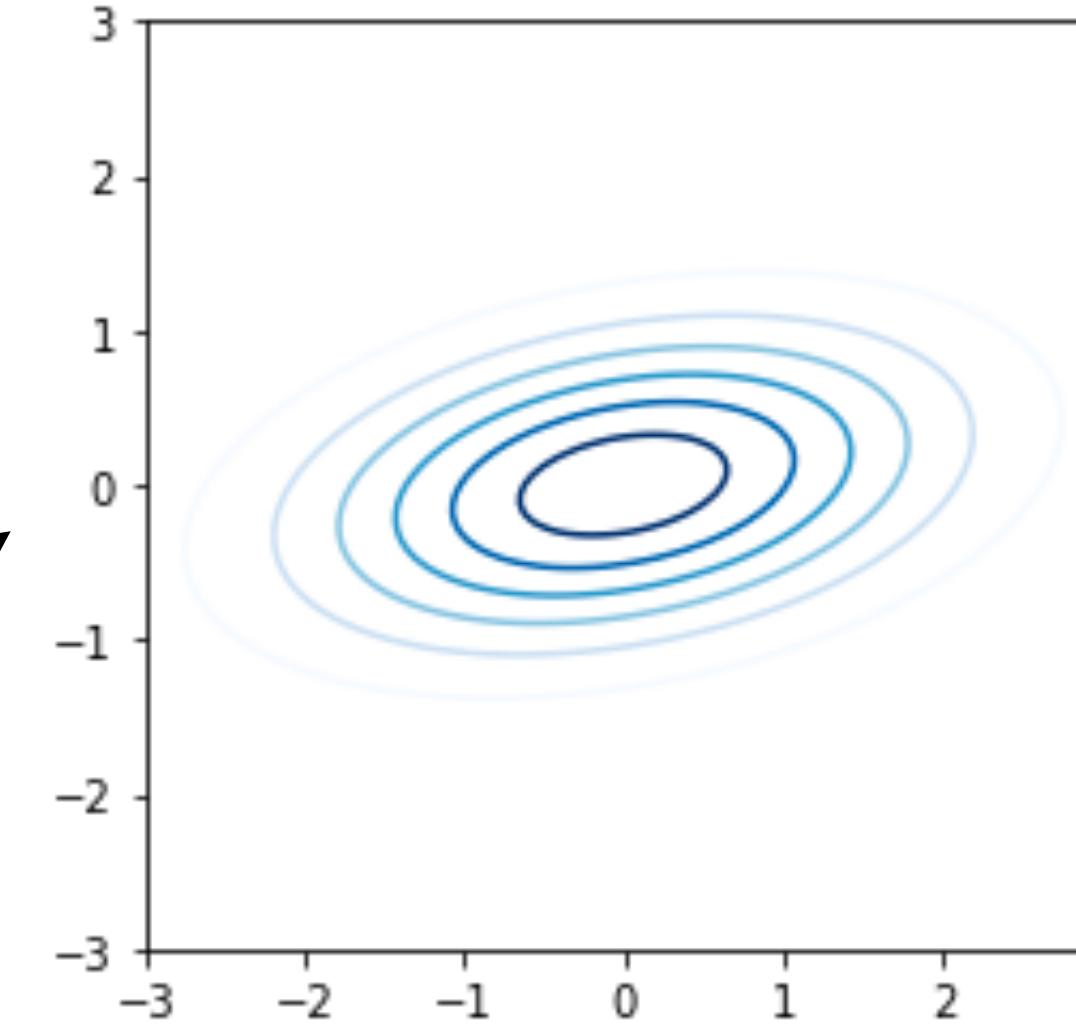
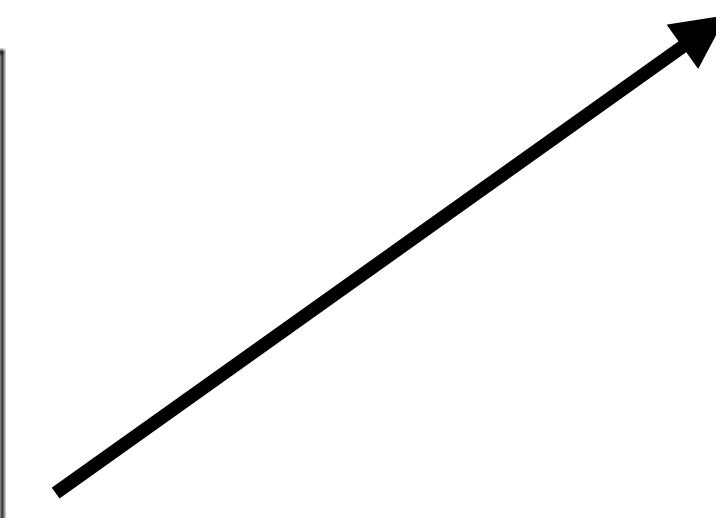
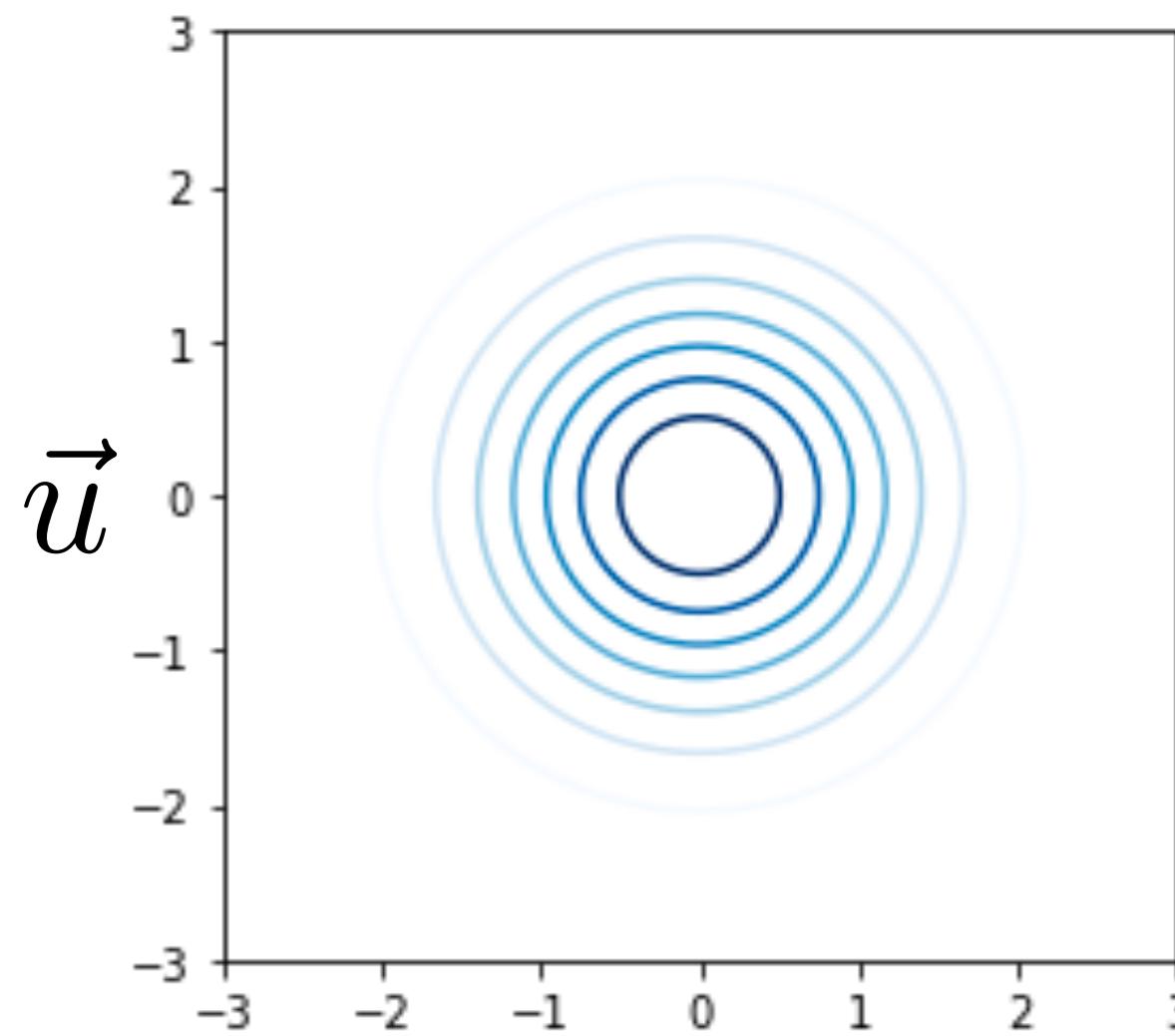
Given $\vec{u} = (u_1, \dots, u_k)$, $u_i \stackrel{iid}{\sim} \mathcal{N}(0, 1)$

$\vec{z} = (z_1, \dots, z_k)$ is M.V.G. iff

there's a matrix R and vector mu
such that

$$\vec{z} = R\vec{u} + \mu$$

Equivalent definition intuition



$$\vec{z} = R\vec{u} + \mu$$

Estimating MVG from data

Maximum likelihood estimate

$$\hat{\mu}_{\text{MLE}}, \hat{\Sigma}_{\text{MLE}} = \arg \max_{\mu, \Sigma} p(\mathcal{D}; \mu, \Sigma)$$

$$= \arg \max_{\mu, \Sigma} \sum_{i=1}^n \log p(x_i; \mu, \Sigma)$$

$$= \arg \min_{\mu, \Sigma} \sum_{i=1}^n \frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) - \frac{n}{2} \log |\Sigma| + C$$

Estimating MVG from data

Maximum likelihood estimate

$$\nabla_{\mu} \left(\sum_{i=1}^n \frac{1}{2} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) - \frac{n}{2} \log |\Sigma| + C \right)$$

$$= \sum_{i=1}^n \nabla_{\mu} \left(\frac{1}{2} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) \right)$$

$$= \sum_{i=1}^n \Sigma^{-1} (x_i - \mu) = 0$$

$$\Rightarrow \sum_{i=1}^n x_i = \sum_{i=1}^n \mu \Rightarrow \mu = \frac{\sum_{i=1}^n x_i}{n}$$

Estimating MVG from data

Maximum likelihood estimate

$$\hat{\mu}_{\text{MLE}} = \frac{\sum_{i=1}^n x_i}{n}$$

$$\hat{\Sigma}_{\text{MLE}} = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu}_{\text{MLE}})(x_i - \hat{\mu}_{\text{MLE}})^T$$

Properties of MVG

Matrix transformations

$$z \sim \mathcal{N}(A\mu_z, Z\Sigma_z A^T)$$

proof:

$$\mu_{Az} = \mathbb{E}[AZ] = A\mathbb{E}[Z] = A\mu_z$$

$$\begin{aligned}\Sigma_{Az} &= \mathbb{E}[(Az - \mathbb{E}[Az])(Az - \mathbb{E}[Az])^T] \\ &= \mathbb{E}[A(z - \mathbb{E}[z])(z - \mathbb{E}[z])^T A^T] \\ &= A\mathbb{E}[(z - \mathbb{E}[z])(z - \mathbb{E}[z])^T] A^T \\ &= A\Sigma_z A^T\end{aligned}$$

Properties of MVG

Marginal distribution

$$\mu_z = \begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}$$

$$\Sigma_z = \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix}$$

$$x \sim \mathcal{N}(\mu_x, \Sigma_{xx})$$

Properties of MVG

Conditional distribution

$x \mid y$

$$\mu_z = \begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix} \quad \Sigma_z = \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix}$$

Properties of MVG

Conditional distribution

$$x \mid y \sim \mathcal{N}$$

$$\mu_z = \begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix} \quad \Sigma_z = \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix}$$

Properties of MVG

Conditional distribution

$$x \mid y \sim \mathcal{N}(\mu_x$$

$$\mu_z = \begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix} \quad \Sigma_z = \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix}$$

Properties of MVG

Conditional distribution

$$x \mid y \sim \mathcal{N}(\mu_x + \Sigma_{xy}\Sigma_{yy}^{-1}(y - \mu_y))$$

$$\mu_z = \begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix} \quad \Sigma_z = \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix}$$

Properties of MVG

Conditional distribution

$$x \mid y \sim \mathcal{N}(\mu_x + \Sigma_{xy}\Sigma_{yy}^{-1}(y - \mu_y), \Sigma_{xx})$$

$$\mu_z = \begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix} \quad \Sigma_z = \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix}$$

Properties of MVG

Conditional distribution

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

$$\mu_z = \begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix} \quad \Sigma_z = \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix}$$

Properties of MVG

Conditional distribution

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

\downarrow

Prior mean

$$\mu_z = \begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix} \quad \Sigma_z = \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix}$$

Properties of MVG

Conditional distribution

$$x \mid y \sim \mathcal{N}(\mu_x + \Sigma_{xy}\Sigma_{yy}^{-1}(y - \mu_y), \underline{\Sigma_{xx} - \Sigma_{xy}\Sigma_{yy}^{-1}\Sigma_{yx}})$$

↓
‘error’ in y

$$\mu_z = \begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix} \quad \Sigma_z = \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix}$$

Properties of MVG

Conditional distribution

$$x \mid y \sim \mathcal{N}(\mu_x + \Sigma_{xy} \underline{\Sigma_{yy}^{-1}(y - \mu_y)}, \Sigma_{xx} - \Sigma_{xy} \Sigma_{yy}^{-1} \Sigma_{yx})$$

↓

‘corrected error’ in y

$$\mu_z = \begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix} \quad \Sigma_z = \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix}$$

Properties of MVG

Conditional distribution

$$x \mid y \sim \mathcal{N}(\mu_x + \underbrace{\Sigma_{xy}\Sigma_{yy}^{-1}(y - \mu_y)}_{\downarrow}, \Sigma_{xx} - \Sigma_{xy}\Sigma_{yy}^{-1}\Sigma_{yx})$$

Implication for x of error in y

$$\mu_z = \begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix} \quad \Sigma_z = \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix}$$

Where to learn more about MVGs

- Previous semester lecture notes
- More rigorous proofs: Michael Jordan's WIP textbook
<http://people.eecs.berkeley.edu/~jordan/prelims/chapter13.pdf>