#### COMP90051

# Workshop Week 10

COPYRIGHT 2017, THE UNIVERSITY OF MELBOURNE

# About the Workshops

- **7** sessions in total
  - **Tue 12:00-13:00 AH211**
  - **Tue 12:00-13:00 AH108 \***
  - **Tue 13:00-14:00 AH210**
  - **Tue 16:15-17:15** AH109
  - **Tue 17:15-18:15 AH236 \***
  - **Tue 18:15-19:15** AH236 \*
  - **Fri** 14:15-15:15 AH211

# About the Workshops

Homepage

https://trevorcohn.github.io/comp90051-2017/workshops

□ Solutions will be released on next Friday (a week later).

COPYRIGHT 2017, THE UNIVERSITY OF MELBOURNE

#### Reminder

**Project** 2

- □ Kaggle competition due on Mon, 09/Oct/17
- □ Worksheet, report, and code due on Wed, 11/Oct/17

**E**xam

**General Fri**, 03/Nov/2017, 8:30am

**3** hours

Royal Exhibition Building

COPYRIGHT 2017, THE UNIVERSITY OF MELBOURNE

# Syllabus

1	Introduction; Probability theory	Probabilistic models; Parameter fitting	
2	Linear regression; Intro to regularization	Logistic regression; Basis expansion	
3	Optimization; Regularization	Perceptron	
4	Backpropagation	CNNs; Auto-encoders	
5	Hard-margin SVMs	Soft-margin SVMs	
6	Kernel methods	Ensemble Learning	
7	Clustering	EM algorithm	
8	Principal component analysis; Multidimensional Scaling	Manifold Learning; Spectral clustering	
9	Bayesian inference (uncertainty, updating)	Bayesian inference (conjugate priors)	$\leftarrow$
10	PGMs, fundamentals	PGMs, independence	
11	Guest lecture (TBC)	PGMs, inference	
12	PGMs, statistical inference	Subject review	

#### Outline

Review the lecture, background knowledge, etc.

- □ MLE, MAP, Bayesian estimates
- Comparison between Bayesian and frequentist
- Likelihood, prior, and posterior
- Conjugate prior and likelihood
  - □ Bayesian linear regression

□ IPython notebook task: Bayesian linear regression

#### MLE, MAP

 $\Box \text{ Training set } \{(\boldsymbol{x}_i, y_i)\}_{i=1}^N, \boldsymbol{X} \text{ for all } \boldsymbol{x}_i, \boldsymbol{y} \text{ for all } y_i \}$ 

 $\square \widehat{\boldsymbol{w}} = \max_{\boldsymbol{w}} \prod_{i=1}^{N} p(y_i | \boldsymbol{x}_i, \boldsymbol{w}) \text{ or } \max_{\boldsymbol{w}} \prod_{i=1}^{N} p(y_i | \boldsymbol{x}_i, \boldsymbol{w}) p(\boldsymbol{w})$  $\square \text{ Prediction for } \boldsymbol{x}^* \text{ is } p(y^* | \boldsymbol{x}^*, \widehat{\boldsymbol{w}})$ 

Choose hyper-parameters / models

on a held-out validation set

□ by cross-validation

on OOB samples (random forest)



$$p(\boldsymbol{w}|\boldsymbol{X}, \boldsymbol{y}) = \frac{p(\boldsymbol{y}|\boldsymbol{X}, \boldsymbol{w})p(\boldsymbol{w})}{p(\boldsymbol{y}|\boldsymbol{X})} \propto \prod_{i=1}^{N} p(y_i|\boldsymbol{x}_i, \boldsymbol{w}) p(\boldsymbol{w})$$

□ Mean estimate E[w], uncertainty  $Var(w) \rightarrow confidence$ □ Prediction for  $x^*$  is  $p(y^*|x^*) = E_{w|X,y} p(y^*|x^*, w)$ 

Choose hyper-parameters / models by comparing p(y|X)

 $\square \text{ main difficulty: get } p(y|X) \text{ or normalize } p(y|X,w)p(w)$ 

 $\Box$  have to approximate the prediction if p(y|X) is intractable

□ some methods can sample from p(w|X, y) without normalizing p(y|X, w)p(w) and then make approximate predictions

# Frequentist and Bayesian

□ Frequentist

- □ find a single parameter vector to best fit the training set
- □ the best parameters are used to make predictions directly

Bayesian

- Germulate the full posterior given the training data
- all the weights are used to make expected predictions
- where each is scaled by its posterior probability

# Bayesian

Advantages

less sensitive to overfitting (expected predictions)

particularly with small training sets

make use of all the data at once

- no need to hold out validation data, or repeatedly train and test
- won't overfit to the held-out set when selecting many parameters

#### Disadvantages

- exact inference is sometimes intractable
- approximate inference may be inefficient and inaccurate
- algorithms are sometimes complex

### Bayesian formula

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

- $\Box p(y|X, w)$  likelihood
- $\Box p(w)$  prior
- $\Box p(y|X)$  marginal likelihood or evidence
- $\Box p(w|X, y)$  posterior

 $\Box p(\boldsymbol{y}|\boldsymbol{X}) = \sum_{\boldsymbol{w}} p(\boldsymbol{y}|\boldsymbol{X}, \boldsymbol{w}) p(\boldsymbol{w}) \text{ or } p(\boldsymbol{y}|\boldsymbol{X}) = \int p(\boldsymbol{y}|\boldsymbol{X}, \boldsymbol{w}) p(\boldsymbol{w}) \, \mathrm{d}\boldsymbol{w}$ 

# Conjugate prior and likelihood

when p(y|X, w)p(w) has the same form as p(w)

□ simplifies the problem of finding the posterior p(w|X, y)□ as needed in Bayesian inference

allows for exact computation of the evidence  $p(\boldsymbol{y}|\boldsymbol{X}) = \frac{p(\boldsymbol{y}|\boldsymbol{X}, \boldsymbol{w})p(\boldsymbol{w})}{p(\boldsymbol{w}|\boldsymbol{X}, \boldsymbol{y})}$ 

# Suite of useful conjugate priors

	likelihood	conjugate prior
n regression	Normal	Normal (for mean)
	Normal	Inverse Gamma (for variance) or Inverse Wishart (covariance)
ificatior	Binomial	Beta
class	Multinomial	Dirichlet
counts	Poisson	Gamma

### **Bayesian Linear Regression (cont)**

- We have two Normal distributions
  \* normal likelihood x normal prior
- Their product is also a Normal distribution
  - \* **conjugate prior:** when product of likelihood x prior results in the same distribution as the prior
  - *evidence* can be computed easily using the normalising constant of the Normal distribution

 $p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \sigma^2) \propto \operatorname{Normal}(\mathbf{w}|\mathbf{0}, \gamma^2 \mathbf{I}_D) \operatorname{Normal}(\mathbf{y}|\mathbf{X}\mathbf{w}, \sigma^2 \mathbf{I}_N)$  $\propto \operatorname{Normal}(\mathbf{w}|\mathbf{w}_N, \mathbf{V}_N)$ 

closed form solution for posterior!

#### **Bayesian Linear Regression (cont)**

 $p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \sigma^2) \propto \operatorname{Normal}(\mathbf{w}|\mathbf{0}, \gamma^2 \mathbf{I}_D) \operatorname{Normal}(\mathbf{y}|\mathbf{X}\mathbf{w}, \sigma^2 \mathbf{I}_N)$  $\propto \operatorname{Normal}(\mathbf{w}|\mathbf{w}_N, \mathbf{V}_N)$ 

where

$$\mathbf{w}_N = \frac{1}{\sigma^2} \mathbf{V}_N \mathbf{X}' \mathbf{y}$$
$$\mathbf{V}_N = \sigma^2 (\mathbf{X}' \mathbf{X} + \frac{\sigma^2}{\gamma^2} \mathbf{I}_D)^{-1}$$

Note that mean (and mode) are the MAP solution from before

Advanced: verify by expressing product of two Normals, gathering exponents together and 'completing the square' to express as squared exponential (i.e., Normal distribution).

#### Outline

Review the lecture, background knowledge, etc.

- □ MLE, MAP, Bayesian estimates
- Comparison between Bayesian and frequentist
- Likelihood, prior, and posterior
- Conjugate prior and likelihood
  - □ Bayesian linear regression

IPython notebook task: Bayesian linear regression