



Language Technologies Institute



# Multimodal Affective Computing

Lecture 9: Discriminative Predictive Modeling

Louis-Philippe Morency Jeffrey Girard

Originally developed with help from Stefan Scherer and Tadas Baltrušaitis

# Outline

- Dynamic Bayesian Network
  - Hidden Markov Models
  - Factorial and coupled HMMs
- Markov Random Fields
  - Unary, binary and clique potentials
  - Factor graph representation
- Multimodal Machine Learning
  - Core Challenges: Representation, Alignment, Fusion, Translation and Co-Learning
- Discriminative Graphical Models
  - Logistic classifier
  - Conditional random fields
  - L1 and L2 regularization
- Evaluation methods and error measures





# **Upcoming Deadlines and Course Schedule**

# Thursday, April 4<sup>th</sup> 4:30pm-6pm

- Midterm presentations
- Sunday, April 7<sup>th</sup> at 11:59pm
  - Midterm report deadline

# Thursday, May 2<sup>nd</sup> 4:30pm-6pm

Final presentations

# Tuesday, May 7<sup>th</sup> at 11:30pm

Final report deadline

\*\*\* No reading assignments for Weeks 12 and 13 \*\*\*



## Midterm Reports – Sunday April 7<sup>th</sup> 11:59pm ET

Main report sections:

- Background and Motivation (1/3 to 1/2 page)
- Literature Review (1/2 to 1 page)
- Data Description and New Annotations (about 1 page)
- Problem Conceptualization (about 1 page)
- Statistical Analysis (1 to 2 pages)
- Next Steps about (1/2 page)
- Appendix: Team Collaboration (about 1/2 page)

# Maximum report length: 7 pages



## Midterm Presentations – Thursday April 4<sup>th</sup>

Main presentation sections:

- Motivation, research problem and dataset
- New annotations
- Problem conceptualization
- Statistical analysis

Presentation instructions:

- Maximum length: 8 minutes
- All teammates should participate
- Followed by questions and feedback forms



## **Upcoming Lectures**

Classes	Tuesday	Thursday
Week 10 3/19 & 3/21 *midterm homework*	<ul> <li>Probabilistic predictive modeling</li> <li>Probabilistic graphical models</li> <li>Bayesian networks and Naïve Bayes classifier</li> <li>Dynamic Bayesian networks and HMMs</li> </ul>	Discussion (probabilistic) • Vaibhav • Vasu Sharma
Week 11 3/26 & 3/28	<ul> <li>Discriminative predictive modeling</li> <li>Markov random fields</li> <li>Factor graph representation</li> <li>Discriminative graphical models</li> </ul>	Discussion (discriminative) <ul> <li>Vaibhav</li> <li>Vasu Sharma</li> </ul>
Week 12 4/02 & 4/04 *midterm report*	<ul> <li>Multimodal deep representations</li> <li>Multimodal joint representations</li> <li>Coordinated representations</li> <li>Temporal representations</li> </ul>	Midterm presentations
Week 13 4/09 & 4/11	<ul> <li>Multimodal alignment and fusion</li> <li>Attention and modality alignment</li> <li>Temporal and multimodal fusion</li> </ul>	NO CLASS
Week 14 4/16 & 4/18	<ul> <li>Multimodal Behavior Generation</li> <li>Guest lecture: Prof. Nakano</li> <li>Generation based on user's attitude</li> <li>Robot and virtual humans</li> </ul>	Discussion (medical) • Jiang Liu • Mahmoud Al Ismail
Week 15 4/23 & 4/25	<ul> <li>Multimodal applications</li> <li>Assessment in the clinical process</li> <li>Biomarkers and behavioral indicators</li> <li>Validation in the medical sciences</li> </ul>	<ul><li>Discussion (educational)</li><li>Mingtong Zhang</li><li>Ankit Shah</li></ul>
Week 16 4/30 & 5/02 *final report*	Final presentations	Final presentations





# Dynamic Bayesian Networks



Language Technologies Institute

# **Dynamic Bayesian Network (DBN)**

- Bayesian network with time-series to represent temporal dependencies.
- Dynamically changing or evolving over time.
- Directed graphical model of stochastic processes.
- Especially aiming at time series modeling.
- Satisfying the Markovian condition: The state of a system at time t depends only on its immediate past state at time t-1.



### **Dynamic Bayesian Network (DBN)**

























# **Factorial HMM**



- Factorial HMM:
  - $h_t$  and  $v_t$  represent two different types of background information, each with its own history
  - Observations x<sub>t</sub> depend on both hidden processes
- Model parameters:
  - $p(v_{t+1}|v_t), p(h_{t+1}|h_t), p(x_t|h_t, v_t)$



# The Boltzmann Zipper



- Both streams have a "memory" ( $h_t$  and  $v_t$ )
- Model parameters:
  - $p(h_{t+1}|h_t), p(x_t|h_t)$
  - $p(v_{t+1}|v_t,h_{t+1}), p(y_t|h_t)$



# The Coupled HMM



- Advantage over Boltzmann Zipper: More flexible, because neither vision nor sound is "privileged" over the other.
  - $p(h_{t+1}|v_t,h_t), p(x_t|h_t)$
  - $p(v_{t+1}|v_t,h_t), p(y_t|h_t)$



### Learning (Dynamic) Bayesian Networks

- Multiple techniques exist to learn the model parameters based on data
  - Maximum likelihood estimator
  - Bayesian estimator, which allows to include prior information
- Python libraries:
  - http://pgmpy.org/
  - http://www.bayespy.org
  - https://pomegranate.readthedocs.io/en/latest/



# Markov Random Field



Language Technologies Institute



**Two Main Types of Graphical Models** 

Bayesian networks (last week)



Markov Models



- Directed acyclic graph
- Conditional dependencies
- Undirected graphical model
- Cyclic dependencies

Langu

### **Markov Random Fields**

$$p(H = \mathbf{h}; \theta) = \frac{\Phi(\mathbf{h}; \theta)}{\sum_{\mathbf{h}'} \Phi(\mathbf{h}'; \theta)}$$

Potential of this variable assignment *h* 

Potential of all possible variable assignments h'

Set of random variables *H* having a Markov property described by undirected graph





### Markov Random Fields – Graphical Model

$$p(H = \mathbf{h}; \theta) = \frac{\Phi(\mathbf{h}; \theta)}{\sum_{\mathbf{h}'} \Phi(\mathbf{h}'; \theta)}$$

$$\Phi(\mathbf{h}; \theta) = \phi_{12}(h_1, h_2; \theta_{12}) \times \phi_{16}(h_1, h_6; \theta_{16}) \times \phi_{26}(h_2, h_6; \theta_{26}) \times \phi_{25}(h_2, h_5; \theta_{25}) \times \phi_{45}(h_4, h_5; \theta_{45}) \times \phi_{34}(h_3, h_4; \theta_{34})$$



 $(h_2)$ 

 $(h_3)$ 

### **Markov Random Fields: Factor Graphs**

$$p(H = h; \theta) = \frac{\Phi(h; \theta)}{\sum_{h'} \Phi(h'; \theta)}$$

$$\Phi(h; \theta) = \phi_{12}(h_1, h_2; \theta_{12}) \times \phi_{16}(h_1, h_6; \theta_{16}) \times \phi_{26}(h_2, h_6; \theta_{26}) \times \phi_{25}(h_2, h_5; \theta_{25}) \times \phi_{45}(h_4, h_5; \theta_{45}) \times \phi_{34}(h_3, h_4; \theta_{34})$$



### Markov Random Fields (Factor Graphs)

$$p(H = h, x; \theta) = \frac{\Phi(h, x; \theta)}{\sum_{h'} \Phi(h', x; \theta)}$$

$$\Phi(h; \theta) = \phi_{12}(h_1, h_2; \theta_{12}) \times$$

$$\phi_{16}(h_1, h_6; \theta_{16}) \times$$

$$\phi_{26}(h_2, h_6; \theta_{26}) \times$$

$$\phi_{16}(h_4, h_5; \theta_{16}) \times$$

$$\phi_{25}(h_2, h_5; \theta_{25}) \times$$

$$\phi_{16}(h_4, h_5; \theta_{45}) \times$$

$$\phi_{34}(h_3, h_4; \theta_{34}) \times$$

$$\psi_1(h_1, x; \theta_1) \times \psi_5(h_5, x; \theta_5) \xrightarrow{\text{Unary potentials}}$$



### Markov Random Fields (Factor Graphs)



### **Markov Random Fields – Clique Factorization**

$$p(H = \mathbf{h}, \mathbf{x}; \theta) = \frac{\Phi(\mathbf{h}, \mathbf{x}; \theta)}{\sum_{\mathbf{h}'} \Phi(\mathbf{h}', \mathbf{x}; \theta)}$$
Clique factorization
$$\Phi(\mathbf{h}; \theta) = \phi_{12}(h_1, h_2; \theta_{12}) \times \phi_{16}(h_1, h_6; \theta_{16}) \times \phi_{26}(h_2, h_6; \theta_{26}) \times \phi_{25}(h_2, h_5; \theta_{25}) \times \phi_{45}(h_4, h_5; \theta_{45}) \times \phi_{34}(h_3, h_4; \theta_{34}) \times \psi_1(h_1; \theta_1) \times \psi_5(h_5; \theta_5)$$

$$\psi_1 = \mathbf{h}, \mathbf{x}; \theta$$

$$\psi_1 = \mathbf{h}, \mathbf{x}; \theta$$

$$\Phi(\mathbf{h}; \theta) = \phi_{12}(h_1, h_2; \theta_{12}) \times \phi_{16}(h_1, h_6; \theta_{16}) \times \phi_{26}(h_2, h_6; \theta_{26}) \times \phi_{25}(h_2, h_5; \theta_{25}) \times \phi_{45}(h_4, h_5; \theta_{45}) \times \phi_{34}(h_3, h_4; \theta_{34}) \times \psi_1(h_1; \theta_1) \times \psi_5(h_5; \theta_5)$$

$$\psi_1 = \mathbf{h}, \mathbf{x}; \theta$$

$$\psi_1 = \mathbf{h}, \mathbf{x}; \theta$$

$$\Phi(\mathbf{h}; \theta) = \phi_{12}(h_1, h_2; \theta_{12}) \times \phi_{26}(h_2, h_6; \theta_{26}) \times \phi_{25}(h_2, h_5; \theta_{25}) \times \phi_{45}(h_3, h_4; \theta_{34}) \times \psi_{16}(h_1; \theta_1) \times \psi_5(h_5; \theta_5)$$

$$\psi_1(h_1; \theta_1) \times \psi_5(h_5; \theta_{55}) \times \phi_{345}(h_3, h_4, h_5; \theta_{345})$$



### **Chain Markov Random Fields (Factor Graphs)**

$$p(H = \mathbf{h}, \mathbf{x}; \theta) = \frac{\Phi(\mathbf{h}, \mathbf{x}; \theta)}{\sum_{\mathbf{h}'} \Phi(\mathbf{h}', \mathbf{x}; \theta)}$$

$$\Phi(\mathbf{h}; \theta) = \phi_{12}(h_1, h_2; \theta_{12}) \times$$

$$\phi_{23}(h_2, h_3; \theta_{23}) \times$$

$$\phi_{34}(h_3, h_4; \theta_{34}) \times$$

$$\psi_1(h_1; \theta_1) \times$$

$$\psi_2(h_2; \theta_2) \times$$

$$\psi_3(h_3; \theta_3) \times$$

$$\psi_4(h_4; \theta_4)$$

$$\psi_4(h_4; \theta_4)$$

$$\psi_4(h_4; \theta_4)$$



### **Example: Markov Random Field – Graphical Model**



Language Technologies Institute

### **Example: Markov Random Field – Factor Graph**



Language Technologies Institute

### **Example: Markov Random Field – Factor Graph**



Language Technologies Institute

### **Example: Markov Random Field – Factor Graph**





# Multimodal Machine Learning: Core Technical Challenges

## **Core Challenges in "Deep" Multimodal ML**

**Representation** 

Alignment

**Fusion** 

**Translation** 

**Co-Learning** 

### Multimodal Machine Learning: A Survey and Taxonomy

By Tadas Baltrusaitis, Chaitanya Ahuja, and Louis-Philippe Morency

https://arxiv.org/abs/1705.09406

✓ 5 core challenges
✓ 37 taxonomic classes
✓ 253 referenced citations

### These challenges are non-exclusive.

**Definition:** Learning how to represent and summarize multimodal data in away that exploits the complementarity and redundancy.







### **Joint Multimodal Representation**





34

# **Joint Multimodal Representations**





Language Technologies Institute

### **Multimodal Vector Space Arithmetic**



[Kiros et al., Unifying Visual-Semantic Embeddings with Multimodal Neural Language Models, 2014]




**Definition:** Learning how to represent and summarize multimodal data in away that exploits the complementarity and redundancy.





#### **Coordinated Representation: Deep CCA**

Learn linear projections that are maximally correlated:



Andrew et al., ICML 2013





#### **Core Challenge 2: Alignment**

**Definition:** Identify the direct relations between (sub)elements from two or more different modalities.



#### Explicit Alignment

The goal is to directly find correspondences between elements of different modalities

#### Implicit Alignment

Uses internally latent alignment of modalities in order to better solve a different problem





#### **Temporal sequence alignment**



Applications:

- Re-aligning asynchronous data

- Finding similar data across modalities (we can estimate the aligned cost)

- Event reconstruction from multiple sources





## **Alignment examples (multimodal)**





#### **Implicit Alignment**



Karpathy et al., Deep Fragment Embeddings for Bidirectional Image Sentence Mapping, https://arxiv.org/pdf/1406.5679.pdf





#### **Core Challenge 3: Fusion**

**Definition:** To join information from two or more modalities to perform a prediction task.



#### 1) Early Fusion



#### 2) Late Fusion





#### **Core Challenge 3: Fusion**

**Definition:** To join information from two or more modalities to perform a prediction task.

B Model-Based (Intermediate) Approaches

- 1) Deep neural networks
- 2) Kernel-based methods
- 3) Graphical models



Multiple kernel learning



Multi-View Hidden CRF



#### **Core Challenge 4: Translation**

**Definition:** Process of changing data from one modality to another, where the translation relationship can often be open-ended or subjective.





#### **Core Challenge 4 – Translation**





Marsella et al., Virtual character performance from speech, SIGGRAPH/Eurographics Symposium on Computer Animation, 2013



**Definition: T**ransfer knowledge between modalities, including their representations and predictive models.







#### **Core Challenge 5: Co-Learning**







## **Taxonomy of Multimodal Research**

#### Representation

- Joint
  - o Neural networks
  - o Graphical models
  - o Sequential
- Coordinated
  - o Similarity
  - o Structured

#### Translation

- Example-based
  - o Retrieval
  - o Combination
- Model-based
  - o Grammar-based

- Encoder-decoder
- Online prediction

## Alignment

- Explicit
  - o Unsupervised
  - Supervised
- Implicit
  - o Graphical models
  - Neural networks

## Fusion

- Model agnostic
  - Early fusion
  - Late fusion
  - Hybrid fusion

- Model-based
  - o Kernel-based
  - o Graphical models
  - Neural networks

## **Co-learning**

- Parallel data
  - Co-training
  - o Transfer learning
- Non-parallel data
  - Zero-shot learning
  - Concept grounding
  - Transfer learning
- Hybrid data
  - Bridging

Tadas Baltrusaitis, Chaitanya Ahuja, and Louis-Philippe Morency, Multimodal Machine Learning: A Survey and Taxonomy



#### Real world tasks tackled by MMML

- Affect recognition
  - Emotion
  - Persuasion
  - Personality traits
- Media description
  - Image captioning
  - Video captioning
  - Visual Question Answering
- Event recognition
  - Action recognition
  - Segmentation
- Multimedia information retrieval
  - Content based/Cross-media















in in black shirt is playing guitar."

construction worker in orange safety vest is working on road."

"two young girls are playing with lego toy."

boy is doing backflip on wakeboard.









(a) answer-phone

(a) get-out-car

(a) fight-person (b) push-up (b) cartwheel









	CHALLENGES				
APPLICATIONS	Representation	TRANSLATION	FUSION	Alignment	CO-LEARNING
Speech Recognition and Synthesis					
Audio-visual Speech Recognition	$\checkmark$		$\checkmark$	$\checkmark$	$\checkmark$
(Visual) Speech Synthesis	$\checkmark$	$\checkmark$			
Event Detection					
Action Classification	$\checkmark$		$\checkmark$		$\checkmark$
Multimedia Event Detection	$\checkmark$		$\checkmark$		$\checkmark$
Emotion and Affect					
Recognition	$\checkmark$		$\checkmark$	$\checkmark$	$\checkmark$
Synthesis	$\checkmark$	$\checkmark$			
Media Description					
Image Description	$\checkmark$	$\checkmark$		$\checkmark$	$\checkmark$
Video Description	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$
Visual Question-Answering	$\checkmark$		$\checkmark$	$\checkmark$	$\checkmark$
Media Summarization	$\checkmark$	$\checkmark$	$\checkmark$		
Multimedia Retrieval					
Cross Modal retrieval	$\checkmark$	$\checkmark$		$\checkmark$	$\checkmark$
Cross Modal hashing	$\checkmark$				$\checkmark$

Tadas Baltrusaitis, Chaitanya Ahuja, and Louis-Philippe Morency, Multimodal Machine Learning: A Survey and Taxonomy





# Discriminative Graphical Models



Language Technologies Institute

#### **Generative versus Discriminative**



Generative or Discriminative?

Answer: It depends on the loss function!

Generative loss function: (joint probability)

Discriminative loss function: (conditional probability)

$$L(\theta) = \sum_{j=1}^{N} P(\mathbf{h}^{(j)}, \mathbf{X}^{(j)}; \theta)$$
$$L(\theta) = \sum_{j=1}^{N} \log P(\mathbf{h}^{(j)} | \mathbf{X}^{(j)}; \theta)$$



#### **Discriminative Model: Logistic classifier**



#### Score function:

$$P(y_t = 1 | \mathbf{x}_t) = \frac{1}{1 + \exp(-\theta \mathbf{x}_t)}$$
$$P(y_t = c | \mathbf{x}_t) = \frac{\exp(\theta_c \mathbf{x}_t)}{\sum_{k=1}^{K} \exp(\theta_k \mathbf{x}_t)}$$

Binary form

Multinomial form



#### **Comparing Linear and Logistic Models**



#### **Discriminative Model: Logistic classifier**



#### **Score function:**

$$P(y_t = c | \mathbf{x}_t) = \frac{\exp(\theta_c \mathbf{x}_t)}{\sum_{k=1}^{K} \theta_c \mathbf{x}_t}$$
 Familiar multinomial form  
$$P(y_t | \mathbf{x}_t) = \frac{1}{Z(\mathbf{x}_t)} \exp\left(\sum_{k=1}^{K} \theta_k f_k(y_t, \mathbf{x}_t)\right)$$
 General form



#### **Discriminative Model: Logistic classifier**









#### **Feature Functions**





#### **Partition Function: Normalizing Constant**





#### **Training and Loss Function**

Label: {0:Dominant, 1:Not-dominant}

Observation vector: [speech-energy, gaze, turn-taking]

$$P(y_t | \boldsymbol{x}_t) = \frac{1}{\mathcal{Z}(\boldsymbol{x}_t)} \exp\left(\sum_{k=1}^K \theta_k f_k(y_t, \boldsymbol{x}_t)\right)$$

Loss function: Conditional log likelihood

$$L(\theta) = \sum_{j=1}^{N} \log P(\mathbf{y}^{(j)} | \mathbf{X}^{(j)}; \theta) - R(\theta)$$



 $y_t$ 

**X**<sub>t</sub>

## Regularization

$$L(\theta) = \sum_{j=1}^{N} \log P(\mathbf{y}^{(j)} | \mathbf{X}^{(j)}; \theta) - R(\theta)$$

- L-2 Norm (Gaussian prior):  $R(\theta) = \lambda \|\theta\|_2$
- L-1 Norm (Laplacian prior):  $R(\theta) = \lambda \|\theta\|_1$





















#### LASSO and ElasticNet



Observation vector: [speech-energy, gaze, turn-taking]

**Lasso loss function:** squared loss with L1 regularization  $L(\theta) = \sum_{j=1}^{N} (y_j - f(x_j; \theta))^2 - \lambda \|\theta\|_1$ 

ElasticNet: squared loss with L1 and L2 regularization

$$L(\theta) = \sum_{j=1}^{N} \left( y_j - f(\mathbf{x}_j; \theta) \right)^2 - \lambda \|\theta\|_1 - \lambda \|\theta\|_2$$



 $\boldsymbol{X}_t$ 

## Conditional Random Fields (CRFs) [McCallum 2001]







#### **Hidden Conditional Random Field**



#### **Learning Multimodal Structure**

#### Modality-private structure

• Internal grouping of observations

#### Modality-shared structure

Interaction and synchrony







#### **Multi-view Latent Variable Discriminative Models**

#### Modality-private structure

Internal grouping of observations

#### Modality-shared structure

Interaction and synchrony



$$p(y|\mathbf{x}^{A}, \mathbf{x}^{V}; \boldsymbol{\theta}) = \sum_{\mathbf{h}^{A}, \mathbf{h}^{V}} p(y, \mathbf{h}^{A}, \mathbf{h}^{V} | \mathbf{x}^{A}, \mathbf{x}^{V}; \boldsymbol{\theta})$$

Approximate inference using loopy-belief



#### **Recap of generative vs discriminative**






## Machine Learning: Evaluation Methods



Language Technologies Institute

### Supervised learning process: two steps

Learning (training): Learn a model using the training data Testing: Test the model using unseen test data to assess the model accuracy





### **Evaluation methods**

- Holdout set: The available data set D is divided into two disjoint subsets,
  - the *training set D<sub>train</sub>* (for learning a model)
  - the test set D<sub>test</sub> (for testing the model)
- Important: training set should not be used in testing and the test set should not be used in learning.
  - Unseen test set provides a unbiased estimate of accuracy.
- The test set is also called the holdout set. (the examples in the original data set D are all labeled with classes.)
- This method is mainly used when the data set *D* is large.
- Unless building person specific models the training and test sets should not contain the same person



### **Evaluation methods (cont...)**

- n-fold cross-validation: The available data is partitioned into *n* equal-size disjoint subsets.
- Use each subset as the test set and combine the rest n-1 subsets as the training set to learn a classifier.
- The procedure is run n times, which give n accuracies.
- The final estimated accuracy of learning is the average of the *n* accuracies.
- 10-fold and 5-fold cross-validations are commonly used.
- This method is used when the available data is not large.



### **Evaluation methods (cont...)**

- Leave-one-out cross-validation: This method is used when the data set is very small.
- It is a special case of cross-validation
- Each fold of the cross validation has only a single test example and all the rest of the data is used in training.
- If the original data has *m* examples, this is *m*fold cross-validation





### **Hyperparameters**

- How do we determine C or  $\gamma$  for SVM training?
- Parameters that we do not learn through optimization are called hyper-parameters
- Need a way to find optimal values for our task
  For some approaches rules of thumb exist
- Need an analytical way to do it
- Common ways
  - Grid search
  - Random search (not as bad as it sounds)





Language Technologies Institute



Error bars: want realistic (conservative) estimates of accuracy

### Take home

- 1. Never touch test data during training/validation
- 2. Never touch test data during training/validation
- 3. Never touch test data during training/validation





# Machine Learning: Measuring Error



Language Technologies Institute

### **Measuring Error**

	Predicted class	
True Class	Yes	No
Yes	TP: True Positive	FN: False Negative
No	FP: False Positive	TN: True Negative

- Error rate = # of errors / # of instances = (FN+FP) / N
- Recall = # of found positives / # of positives

= TP / (TP+FN) = sensitivity = hit rate

- Precision = # of found positives / # of found
  - = TP / (TP+FP)
- Specificity = TN / (TN+FP)
- False alarm rate = FP / (FP+TN) = 1 Specificity



### **F**<sub>1</sub>-value (also called **F**<sub>1</sub>-score)

 It is hard to compare two classifiers using two measures. F<sub>1</sub> score combines precision and recall into one measure

• 
$$F_1 = \frac{2 \cdot p \cdot r}{p + r}$$

•  $F_1$  - score is the harmonic mean of precision and recall

• 
$$F_1 = \frac{2}{\frac{1}{p} + \frac{1}{r}}$$

- The harmonic mean of two numbers tends to be closer to the smaller of the two
- Preferred over accuracy when data is unbalanced
  - Why?





#### **Receiver Operating Characteristic (ROC) Curve**





### **AUC for ROC curves**





### **Evaluation of regression**

- Root Mean Square Error
  - $\sqrt{\sum_i (y_i x_i)^2}$
  - Not easily interpretable
- Correlation trend prediction in a way
  - Nice interpretation: 0 no relationship, 1 perfect relationship

• 
$$\rho = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{(n-1)\sigma_x \sigma_y}$$

- Concordance Correlation Coefficient (CCC)
  - A method to combine both

• 
$$\rho_c = \frac{2\rho\sigma_x\sigma_y}{\sigma_x^2 + \sigma_y^2 + (\mu_x - \mu_y)^2}$$
,  $\rho$  – correlation coefficient

Has nice interpretability as well



### Take home

- Error measure selection is not straightforward
  - Pick the right one for your problem
  - F1, AUC, Accuracy, RMSE, CCC
- Make sure the same measure is used for validation and testing
  - Otherwise you might be learning suboptimal models
- Wrong error measure can hide both bad and good results

