# Bioinformatics:
# Predicting Networks

Mark Gerstein, Yale University
gersteinlab.org/courses/452
(last edit in Spring '18)

# **Origin of Biological Networks**

# Origin of Networks

- Protein-protein interactions
  - ◊ Phosphorylation networks

- Metabolic Networks

- Regulatory networks
  - ◊ from Chip-Seq (see next slide)

- "Squared" scale
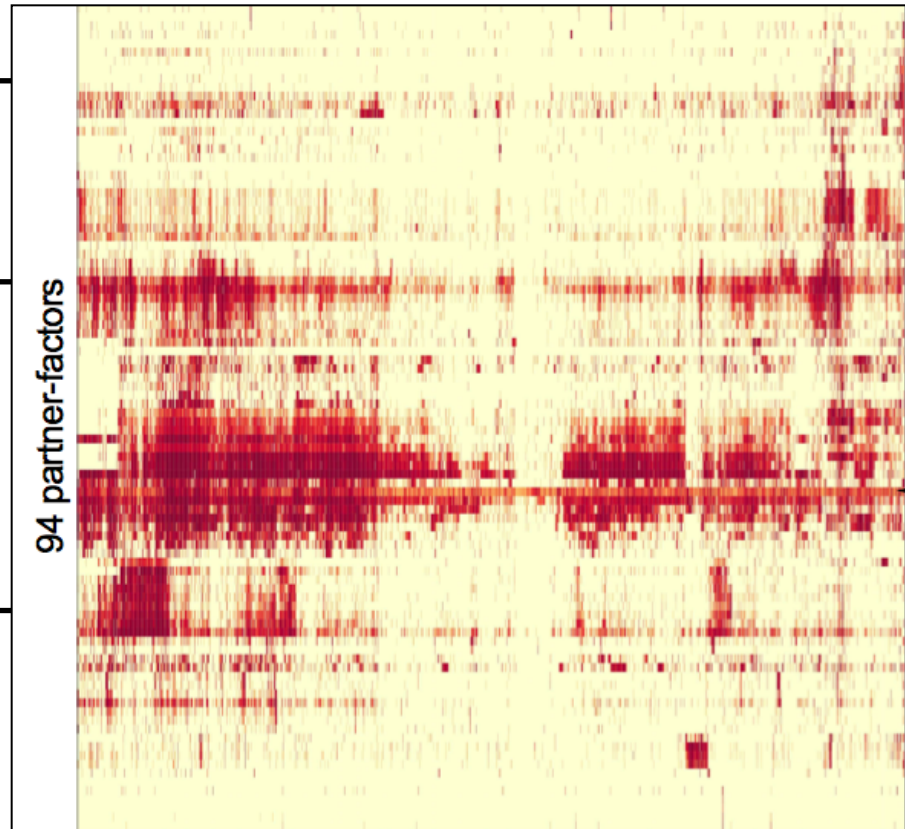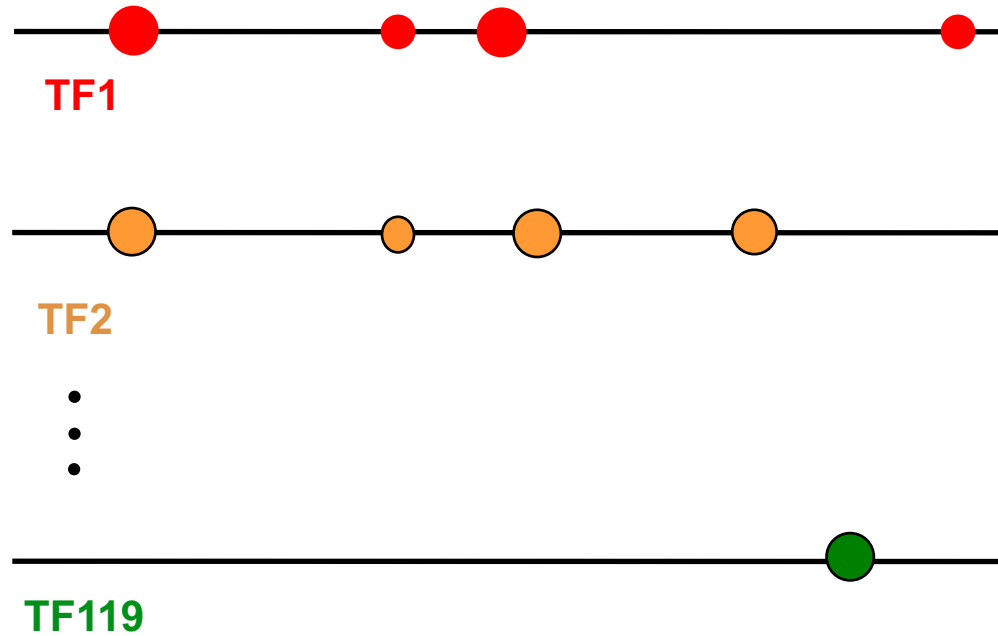  - ◊ 6K genes in yeast but ~18M potential interactions (6000 chose 2 pairs of interactions)

# Data Flow: Chip-seq expts. to co-associating peaks

**119 TFs** from 458 ChIP-Seq experiments (2 Tb tot.)

Signal Tracks

**7M Peaks** from Uniform Peak Calling

TF1

TF2

TF119



94 partner-factors

2785 GATA1 (focus-factor) peak locations

[ Gerstein et al. Nature (in press, '12) ]

# Data Flow: peaks to proximal & distal networks

**Peak Calling**

**Assigning TF binding sites to targets**

~500K
Edges

**Filtering high confidence edges & distal regulation**

Based on stat. model combining
signal strength & location relative to typical binding

~26K
Edges

**TF**

**Potential Distal Edge**

**TF**

**Strong Proximal Edge**

# **Predicting Networks via Bayesian Integration: Problem Motivation**

# RNA polymerase II: Structure

## Which subunits interact?
### Based on Binding experiments



**Source: Edwards et al., 2002, *Trends in Genetics***

## Compare with Gold Std. Structure



**Source: Cramer et al., 2000, *Science*, 288:632-633**

**Subunits**  1 1 1 1 1 1 1 1 1 2 2 2 2 2 2 2 2 3 3 3 3 3 3 3 5 5 5 5 5 5 6 6 6 6 6 8 8 8 8 9 9 9 10 10 12
**Subunits**  2 3 5 6 8 9 10 11 12 3 5 6 8 9 10 11 12 5 6 8 9 10 11 12 6 8 9 10 11 12 8 9 10 11 12 9 10 11 12 10 11 12 11 11 12

**Pull-down 1**   1 1 0 1 0 1 0   1 1 0 1 0 1 0   1 1 1 0 1 1   1 1 0 1 0   0 0 0 0   0 1 0   0 0   0

**Pull-down 2**   1 1 1 1 0 1 0   1 1 0 1 0 1 0   1 1 1 0 1 1   0 1 0 1 0   0 0 0 0   0 0 0   0 0   0

**Pull-down 3**   1           1           1 0 1 0 0 1 0

**Cross-linking**  1 1 1 1   0 1 1 1 1 1 0   1 1 1       1 1   1 0         1

**Far Western 1**   1 1           1 1         1 0 0   0 0 0 0 1   0 0 0

**Far Western 2**     1 1   1 1 1     1 1   1 1 1   0 0   0 1 0 0 0   0 0 0   0 0 0   0 0 0         0 0 0

**Far Western 3**               1 0 0   0 1 0

# Interaction experiments
***before* structure was known**

# Gold-Standard Positives

**Subunits**

| 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 5 | 5 | 5 | 5 | 5 | 5 | 6 | 6 | 6 | 6 | 6 | 8 | 8 | 8 | 8 | 9 | 9 | 9 | 10 | 10 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|----|----|----|
| 2 | 3 | 5 | 6 | 8 | 9 | 10 | 11 | 12 | 3 | 5 | 6 | 8 | 9 | 10 | 11 | 12 | 5 | 6 | 8 | 9 | 10 | 11 | 12 | 6 | 8 | 9 | 10 | 11 | 12 | 8 | 9 | 10 | 11 | 12 | 9 | 10 | 11 | 12 | 9 | 10 | 11 | 12 | 10 | 11 | 11 | 11 | 12 |

**Subunits**

Gold-Standard Positive (GSTD+): 13

# Gold-Standard Negatives

**Subunits**

| 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 5 | 5 | 5 | 5 | 5 | 5 | 6 | 6 | 6 | 6 | 6 | 8 | 8 | 8 | 8 | 9 | 9 | 9 | 10 | 10 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

**Subunits**

| 2 | 3 | 5 | 6 | 8 | 9 | 10 | 11 | 12 | 3 | 5 | 6 | 8 | 9 | 10 | 11 | 12 | 5 | 6 | 8 | 9 | 10 | 11 | 12 | 6 | 8 | 9 | 10 | 11 | 12 | 8 | 9 | 10 | 11 | 12 | 9 | 10 | 11 | 12 | 10 | 11 | 12 | 11 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

Gold-Standard Negative (GSTD-): 32

| Subunits | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 5 | 5 | 5 | 5 | 5 | 5 | 6 | 6 | 6 | 6 | 6 | 8 | 8 | 8 | 8 | 9 | 9 | 9 | 10 | 10 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Subunits | 2 | 3 | 5 | 6 | 8 | 9 | 10 | 11 | 12 | 3 | 5 | 6 | 8 | 9 | 10 | 11 | 12 | 5 | 6 | 8 | 9 | 10 | 11 | 12 | 6 | 8 | 9 | 10 | 11 | 12 | 8 | 9 | 10 | 11 | 12 | 9 | 10 | 11 | 12 | 10 | 11 | 12 | 11 | 11 | 12 |

Gold-Standard Positive (GSTD+): 13

Gold-Standard Negative (GSTD-): 32

# Assess Quality and Coverage of PPints

# Data integration: RNA polymerase II

**Subunit A**

1 1 1 1 1 1 1 1 1 2 2 2 2 2 2 2 2 3 3 3 3 3 3 3 5 5 5 5 5 5 6 6 6 6 6 8 8 8 8 9 9 9 10 10 11

**Subunit B**

2 3 5 6 8 9 10 11 12 3 5 6 8 9 10 11 12 5 6 8 9 10 11 12 6 8 9 10 11 12 8 9 10 11 12 9 10 11 12 10 11 12 11 12 12

**structural contact**

1 0 1 1 1 1 0 1 0 1 0 0 0 1 1 0 1 0 0 0 0 1 1 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0

**Far western**

1 1   1 1   1 0 0   0 0 0 1   0 0 0

**Cross-linking**

1 1 1 1 1   0 1 1 1 1 0   1 1 1   1 1   1 0   1

**Far western**

1 1   1 1 1   1 1   1 1 1   0 0   0 1 0 0 0   0 0 0 0   0 0 0   0 0 0   0 0 0

**Pull-down**

1 1 0 1 0 1 0   1 1 0 1 0 1 0   1 1 1 0 1 1   1 1 0 1 0   0 0 0 0   0 1 0   0 0   0

**Pull-down**

1 1 1 1 0 1 0   1 1 0 1 0 1 0   1 1 1 0 1 1   0 1 0 1 0   0 0 0 0   0 0 0   0 0   0

**Pull-down**

1   1   1 0 1 0 0 1 0

**Far western**

1 0 0   0 1 0

■ = false

■ = true

# Data integration: RNA polymerase II



Subunit A

Subunit B

structural contact

Far western

Cross-linking

Far western

Pull-down

Pull-down

Pull-down

Far western

= false

= true

Union

# Data integration: RNA polymerase II

**Subunit A**    1 1 1 1 1 1 1 1 1 2 2 2 2 2 2 2 2 3 3 3 3 3 3 3 5 5 5 5 5 5 6 6 6 6 6 8 8 8 8 9 9 9 10 10 11

**Subunit B**    2 3 5 6 8 9 10 11 12 3 5 6 8 9 10 11 12 5 6 8 9 10 11 12 6 8 9 10 11 12 8 9 10 11 12 9 10 11 12 10 11 12 11 12 12
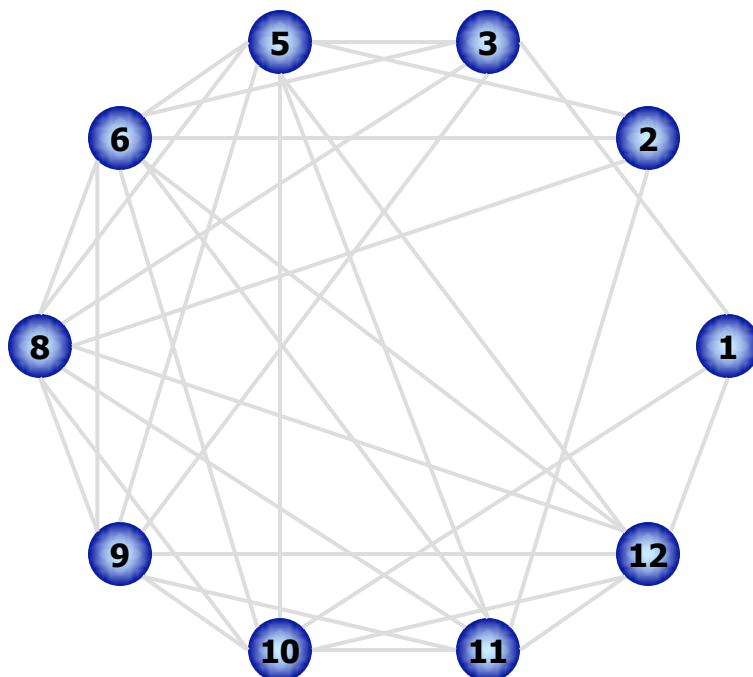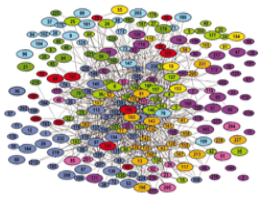
**structural contact**    1 0 1 1 1 1 0 1 0 1 0 0 0 1 1 0 1 0 0 0 0 1 1 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0

**Far western**

**Cross-linking**

**Far western**

**Pull-down**

**Pull-down**

**Pull-down**

**Far western**

**Majority**    1 1 1 1 0 1 0 1 1 1 0 1 0 1 0 1 1 0 0 0 0 1 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0

**Intersection**    1 1 0 1 0 0 0 1 1 1 0 0 0 1 0 1 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0

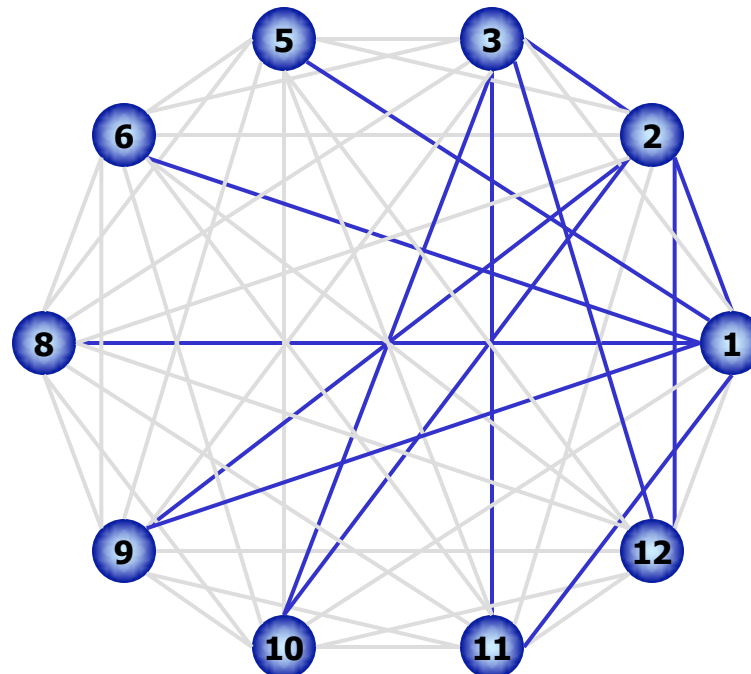**Union**    1 1 1 1 0 1 1 1 1 1 1 0 1 1 1 1 1 0 1 1 0 1 1 0 1 1 0 0 1 0 0 1 0 0 0 0 0 0 1 0 0 0 0 0 0
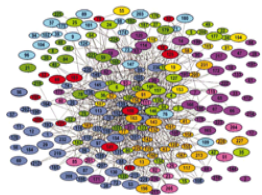
# Data integration: RNA polymerase II

**Subunit A**
1 1 1 1 1 1 1 1 1 2 2 2 2 2 2 2 2 3 3 3 3 3 3 3 5 5 5 5 5 5 6 6 6 6 6 8 8 8 8 9 9 9 10 10 11

**Subunit B**
2 3 5 6 8 9 10 11 12 3 5 6 8 9 10 11 12 5 6 8 9 10 11 12 6 8 9 10 11 12 8 9 10 11 12 9 10 11 12 10 11 12 11 12 12
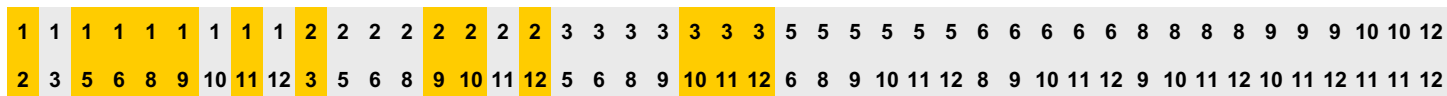
**structural contact**
1 0 1 1 1 1 0 1 0 1 0 0 0 1 1 0 1 0 0 0 0 1 1 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0

**Far western**
1 1    1 1    1 0 0   0 0 0 0 1   0 0 0

**Cross-linking**
1 1 1 1 1   0 1 1 1 1 0   1 1 1    1 1   1 0     1

**Far western**
  1 1   1 1    1 1   1 1 1   0 0   0 1 0 0 0   0 0 0 0   0 0 0   0 0 0    0 0 0

**Pull-down**
1 1 0 1 0 1 0   1 1 0 1 0 1 0   1 1 1 0 1 1   1 1 0 1 0   0 0 0 0   0 1 0    0 0    0

**Pull-down**
1 1 1 1 0 1 0   1 1 0 1 0 1 0   1 1 1 0 1 1   0 1 0 1 0   0 0 0 0   0 0 0    0 0    0

**Pull-down**
1        1       1 0 1 0 0 1 0

**Far western**
1 0 0   0 1 0

(Cross validate)

Integrate using naive Bayes classifier

**Combined (Bayesian)**
0 1 1 1 1 0 0 0 1 1 1 0 0 0 1 0 1 0 0 0 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
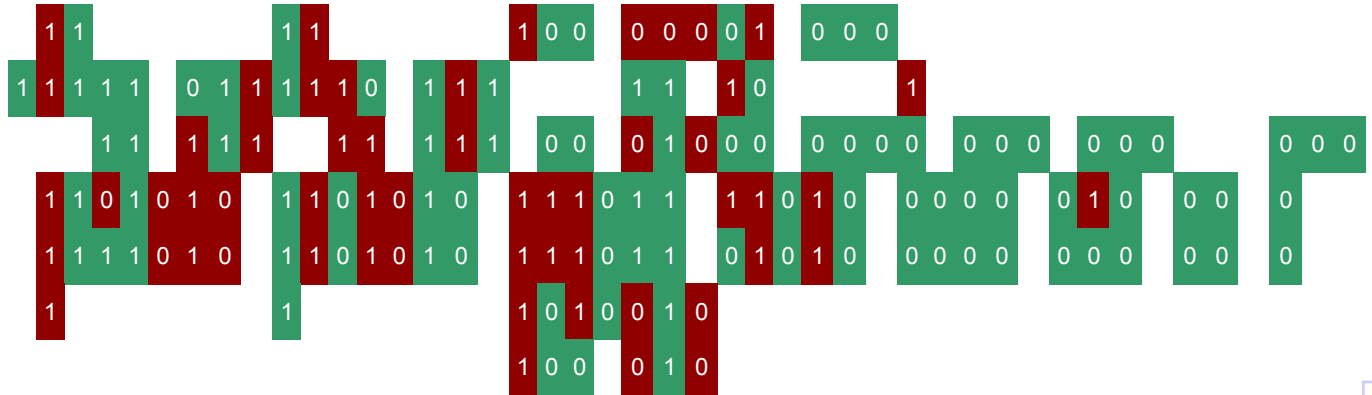
**Majority**
1 1 1 1 1 0 1 0 1 1 1 0 1 0 1 0 1 1 0 0 0 0 1 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0

**Intersection**
1 1 0 1 0 0 0 1 1 1 0 0 0 1 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0

**Union**
1 1 1 1 1 0 1 1 1 1 1 1 0 1 1 1 1 1 0 1 1 0 1 1 0 1 0 0 1 0 0 0 0 0 0 1 0 0 0 0 0 0 0 0 0

# Weighted Voting: the Likelihood Ratio



**structural contact**   1 0 1 1 1 1 0 1 0 1 0 0 0 1 1 0 1 0 0 0 1 1 1 0 0 0 0

**Far western**
**Far western (dup)**
**Cross-linking**
**Far western**
**Pull-down**
**Pull-down**
**Pull-down**
**Far western**

**Combined**   0 1 1 1 1 0 0 0 1 1 1 0 0 0 1 0 1 0 0 0 1 0 0 0 0 0

Maj. Vote: $0$ = round(avg( $0 + 0 + 0 + 1 + 1 + 0 + 0$ )

With weights: **likelihood ratio L** = $L_1 + L_2 + L_3 \ldots$

# **<u>Predicting Networks via Bayesian Integration: Intuition & Formalism</u>**

Derived from "perceptron model"
R = <w,f> + b

**Simple Vote:** $R = f_1 + f_2 + f_3 + ... + f_n$    With $f$ = 1 or -1

**If** $\begin{cases} R>0; & I \text{ Interact} \\ R<0; & \sim I \text{ No interaction} \end{cases}$

**Modify with feature weight:**

$$R = w_1 f_1 + w_2 f_2 + w_3 f_3 + ... + w_n f_n = \vec{w} \cdot \vec{f}$$

If has prior knowledge $w_0$

$$R = \vec{w} \cdot \vec{f} + w_0$$

# Classification by Voting

$$R = \vec{w} \cdot \vec{f} + w_0$$

$$w_1 = \log \frac{P(f_1 = 1 \mid I)}{P(f_1 = 1 \mid \sim I)}$$

$$= \log \frac{TP / P}{FR / N}$$

$$w_0 = \log \frac{P}{N}$$ (Estimated from Golden Standard)

# Bayes Rule

$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)}$$

## Which is shorthand for:

$$P(Y = y_i | X = x_j) = \frac{P(X = x_j | Y = y_i)P(Y = y_i)}{P(X = x_j)}$$

[From Mitchell, Machine Learning]

**Thus**

$$P(I \mid f_1, f_2, f_3, \ldots) = \frac{P(f_1, f_2, f_3, \ldots \mid I)P(I)}{P(f_1, f_2, f_3, \ldots)}$$

Assume Naïve Bayes (independent)

$$= \frac{P(f_1 \mid I)P(f_2 \mid I)P(f_3 \mid I)\ldots P(I)}{P(f_1, f_2, f_3, \ldots)}$$

$$P(\sim I \mid f_1, f_2, f_3, \ldots) = \frac{P(f_1, f_2, f_3, \ldots \mid \sim I)P(\sim I)}{P(f_1, f_2, f_3, \ldots)}$$

$$= \frac{P(f_1 \mid \sim I)P(f_2 \mid \sim I)P(f_3 \mid \sim I)\ldots P(\sim I)}{P(f_1, f_2, f_3, \ldots)}$$

$$\log\left(\frac{P(I \mid f_1, f_2, f_3, \ldots)}{P(\sim I \mid f_1, f_2, f_3, \ldots)}\right) = \log\left(\frac{P(f_1 \mid I)}{P(f_1 \mid \sim I)} \frac{P(f_2 \mid I)}{P(f_2 \mid \sim I)} \frac{P(f_3 \mid I)}{P(f_3 \mid \sim I)} \ldots \frac{P(I)}{P(\sim I)}\right)$$

$$= \log\frac{TPR_1}{FPR_1} + \log\frac{TPR_2}{FPR_2} + \log\frac{TPR_3}{FPR_3} + .. + \log\frac{P}{N}$$

# More Bayes Rule

$$\log\left(\frac{P(I \mid f_1, f_2, f_3, \ldots)}{P(\sim I \mid f_1, f_2, f_3, \ldots)}\right) = \log\frac{TPR_1}{FPR_1} + \log\frac{TPR_2}{FPR_2} + \log\frac{TPR_3}{FPR_3} + \ldots + \log\frac{P}{N}$$

$$w_1 \qquad\qquad w_2 \qquad\qquad w_3 \qquad\qquad w_0$$

# More Bayes Rule

# **Predicting Networks via Bayesian Integration: Worked Examples**

# Likelihood Ratios

**Subunits**

| 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 5 | 5 | 5 | 5 | 5 | 5 | 6 | 6 | 6 | 6 | 6 | 8 | 8 | 8 | 8 | 9 | 9 | 9 | 10 | 10 | 12 |

**Subunits**

| 2 | 3 | 5 | 6 | 8 | 9 | 10 | 11 | 12 | 3 | 5 | 6 | 8 | 9 | 10 | 11 | 12 | 5 | 6 | 8 | 9 | 10 | 11 | 12 | 6 | 8 | 9 | 10 | 11 | 12 | 8 | 9 | 10 | 11 | 12 | 9 | 10 | 11 | 12 | 10 | 11 | 12 | 11 | 11 | 12 |

**Pull-down 1**

| 1 | 1 | 0 | 1 | 0 | 1 | 0 | | 1 | 1 | 0 | 1 | 0 | 1 | 0 | | 1 | 1 | 1 | 0 | 1 | 1 | | 1 | 1 | 0 | 1 | 0 | | 0 | 0 | 0 | 0 | | 0 | 1 | 0 | | 0 | 0 | | 0 |

$$L_1 = \frac{p\left(x_1 \mid GSTD+\right)}{p\left(x_1 \mid GSTD-\right)}$$

$$L_0 = \frac{p\left(x_0 \mid GSTD+\right)}{p\left(x_0 \mid GSTD-\right)}$$

Likelihood Ratio
for Feature $f$:

$$L_f \equiv \frac{p\left(x_f \mid GSTD+\right)}{p\left(x_f \mid GSTD-\right)}$$

**GSTD+**
**GSTD-**
**True**
**False**

**Subunits**

| 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 5 | 5 | 5 | 5 | 5 | 5 | 6 | 6 | 6 | 6 | 6 | 8 | 8 | 8 | 8 | 9 | 9 | 9 | 10 | 10 | 12 |

**Subunits**

| 2 | 3 | 5 | 6 | 8 | 9 | 10 | 11 | 12 | 3 | 5 | 6 | 8 | 9 | 10 | 11 | 12 | 5 | 6 | 8 | 9 | 10 | 11 | 12 | 6 | 8 | 9 | 10 | 11 | 12 | 8 | 9 | 10 | 11 | 12 | 9 | 10 | 11 | 12 | 10 | 11 | 12 | 11 | 11 | 12 |

**Pull-down 1**

| 1 | 0 | 1 | 0 | | 0 | | 1 | | | 0 | 1 | | | | 1 | 1 | |

$$L_1 = \frac{p\left(x_1 \mid GSTD+\right)}{p\left(x_1 \mid GSTD-\right)} = \frac{6/13}{\underline{\phantom{xxxxx}}}$$

$$L_0 = \frac{p\left(x_0 \mid GSTD+\right)}{p\left(x_0 \mid GSTD-\right)} = \frac{4/13}{\underline{\phantom{xxxxx}}}$$

■ **GSTD+**
■ **GSTD-**
■ **True**
■ **False**

**Subunits**

| 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 5 | 5 | 5 | 5 | 5 | 5 | 6 | 6 | 6 | 6 | 6 | 8 | 8 | 8 | 8 | 9 | 9 | 9 | 10 | 10 | 12 |

**Subunits**

| 2 | 3 | 5 | 6 | 8 | 9 | 10 | 11 | 12 | 3 | 5 | 6 | 8 | 9 | 10 | 11 | 12 | 5 | 6 | 8 | 9 | 10 | 11 | 12 | 6 | 8 | 9 | 10 | 11 | 12 | 8 | 9 | 10 | 11 | 12 | 9 | 10 | 11 | 12 | 10 | 11 | 12 | 11 | 11 | 12 |

**Pull-down 1**

1 1 0 1 0 1 0    1 1 0 1 0 1 0    1 1 1 0 1 1    1 1 0 1 0    0 0 0 0    0 1 0    0 0    0

$$L_1 = \frac{p(x_1 \mid GSTD+)}{p(x_1 \mid GSTD-)} = \frac{6/13}{11/32} = 1.34$$

$$L_0 = \frac{p(x_0 \mid GSTD+)}{p(x_0 \mid GSTD-)} = \frac{4/13}{14/32} = 0.70$$

**GSTD+**

**GSTD-**

**True**

**False**

# Calculating Likelihood Ratios

**Subunits**

The slide contains a heatmap table with subunit pairings and binary values for each experimental method, followed by likelihood ratio calculations.

**Pull-down 1** $L1 = (6/13) / (11/32) = 1.34 \quad L0 = (4/13) / (14/32) = 0.70$

**Pull-down 2** $L1 = (7/13) / (9/32) = 1.91 \quad L0 = (2/13) / (16/32) = 0.31$

**Pull-down 3** $L1 = (2/13) / (3/32) = 1.64 \quad L0 = (2/13) / (2/32) = 2.46$

**Cross-linking** $L1 = (10/13) / (7/32) = 3.52 \quad L0 = (0/13) / (3/32) = 0$

**Far Western 1** $L1 = (2/13) / (4/32) = 1.23 \quad L0 = (3/13) / (6/32) = 1.23$

**Far Western 2** $L1 = (6/13) / (5/32) = 2.95 \quad L0 = (2/13) / (17/32) = 0.29$

**Far Western 3** $L1 = (1/13) / (1/32) = 2.46 \quad L0 = (2/13) / (2/32) = 2.46$

Legend:
- GSTD+
- GSTD-
- True
- False

# Data Integration: ROC-Curve

**Subunits**

| 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 5 | 5 | 5 | 5 | 5 | 5 | 6 | 6 | 6 | 6 | 6 | 8 | 8 | 8 | 8 | 9 | 9 | 9 | 10 | 10 | 12 |
| 2 | 3 | 5 | 6 | 8 | 9 | 10 | 11 | 12 | 3 | 5 | 6 | 8 | 9 | 10 | 11 | 12 | 5 | 6 | 8 | 9 | 10 | 11 | 12 | 6 | 8 | 9 | 10 | 11 | 12 | 8 | 9 | 10 | 11 | 12 | 9 | 10 | 11 | 12 | 10 | 11 | 12 | 11 | 11 | 12 |

**Pull-down 1** 1 1 0 1 0 1 0   1 1 0 1 0 1 0   1 1 1 0 1 1   1 1 0 1 0   0 0 0 0   0 1 0   0 0 0   0

**Pull-down 2** 1 1 1 1 0 1 0   1 1 0 1 0 1 0   1 1 1 0 1 1   0 1 0 1 0   0 0 0 0   0 0 0   0 0 0   0

**Pull-down 3** 1   1   1 0 1 0 0 1 0

**Cross-linking** 1 1 1 1 1   0 1 1 1 1 0   1 1 1   1 1   1 0   1

**Far Western 1** 1 1   1 1   1 0 0   0 0 0 1   0 0 0

**Far Western 2** 1 1   1 1 1   1 1   1 1 1   0 0   0 1 0 0   0 0 0 0   0 0 0   0 0 0   0 0 0

**Far Western 3** 1 0 0   0 1 0

Values: 3.52 18.2 11.1 13.9 26.6 0.22 0 2.25 10.4 18.2 11.1 2.25 0 0.22 26.6 2.25 10.4 12.7 5.52 3.68 0.53 19.5 132 2.16 0.52 0 0.22 0.91 0.08 0.36 0.22 0.22 0.06 0.06 0.29 0.22 0.12 0.06 0.29 0.22 0.22 0 0.06 0.29 0.29

**Combined (Bayes)** 1 1 1 1 1 0 0 0 1 1 1 0 0 0 1 0 1 1 0 0 0 1 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0

$$L(f_1, ..., f_n) = L(f_1)...L(f_n)$$

**GSTD+**
**GSTD-**
**True**
**False**

"Weighted Voting"

# Data Integration: ROC Curve

| Subunits | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 5 | 5 | 5 | 5 | 5 | 5 | 6 | 6 | 6 | 6 | 6 | 8 | 8 | 8 | 8 | 9 | 9 | 9 | 10 | 10 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Subunits | 2 | 3 | 5 | 6 | 8 | 9 | 10 | 11 | 12 | 3 | 5 | 6 | 8 | 9 | 10 | 11 | 12 | 5 | 6 | 8 | 9 | 10 | 11 | 12 | 6 | 8 | 9 | 10 | 11 | 12 | 8 | 9 | 10 | 11 | 12 | 9 | 10 | 11 | 12 | 10 | 11 | 12 | 11 | 11 | 12 |

| | 3.52 | 18.2 | 11.1 | 13.9 | 26.6 | 0.22 | 0 | 2.25 | 10.4 | 18.2 | 11.1 | 2.25 | 0 | 0.22 | 26.6 | 2.25 | 10.4 | 12.7 | 5.52 | 3.68 | 0.53 | 19.5 | 132 | 2.16 | 0.52 | 0 | 0.22 | 0.91 | 0.08 | 0.36 | 0.22 | 0.22 | 0.06 | 0.06 | 0.29 | 0.22 | 0.12 | 0.06 | 0.29 | 0.22 | 0.22 | 0 | 0.06 | 0.29 | 0.29 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Combined (Bayes) | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Majority | 1 | 1 | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Intersection | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Union | 1 | 1 | 1 | 1 | 1 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 1 | 1 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

GSTD+
GSTD-
True
False

ROC Curves of RNA Polymerase II

TPR=TP/P=Sensitivity

FPR=FP/N=1-Specificity

○ Bayesian Integration
● Majority
● Intersection
● Union