

Model Selection in Machine Learning

+

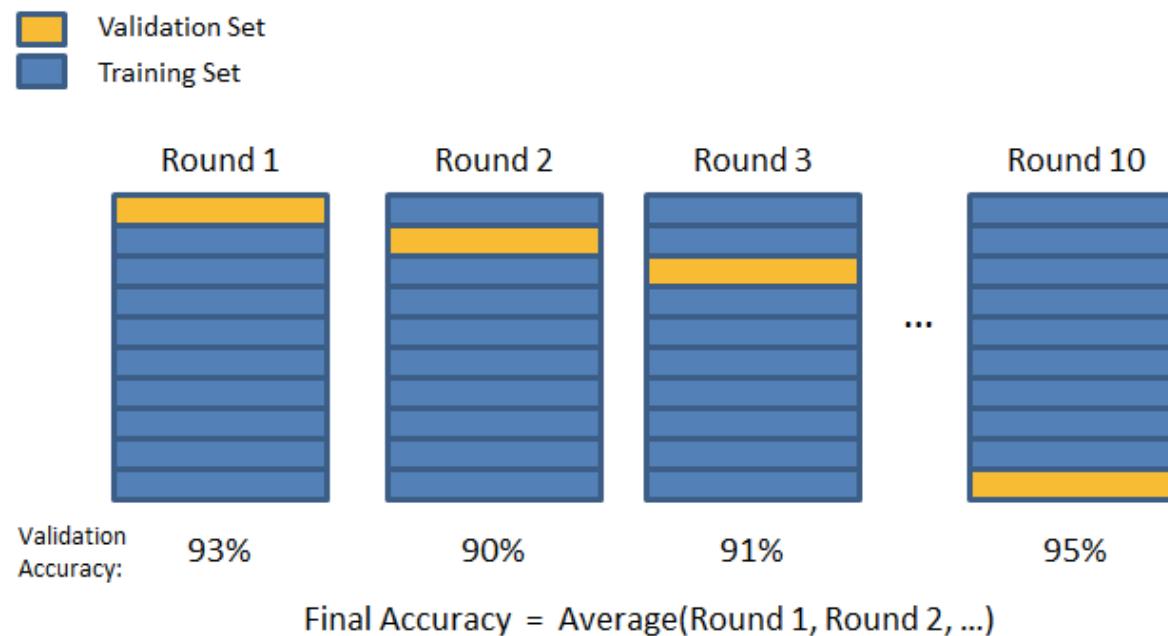
Predicting Gene Expression from ChIP-Seq signals



Review and warm-up questions

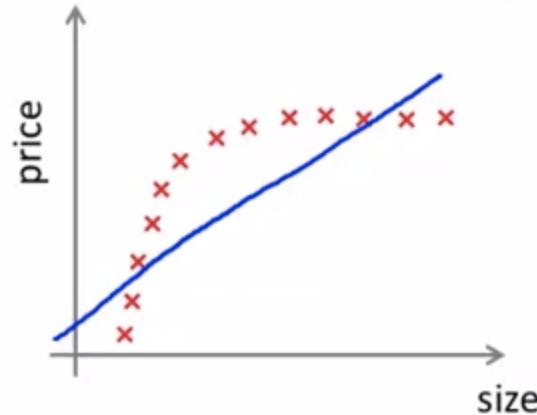
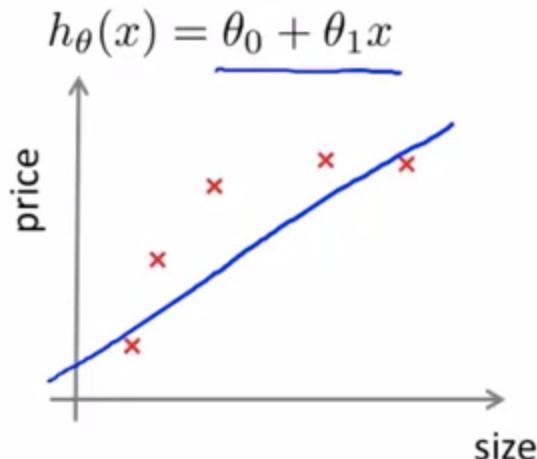
Training vs. Cross-Validation

- Fit model to example data points
- Evaluate model on *separate* set of data points



Bias vs. Variance

Model too simple

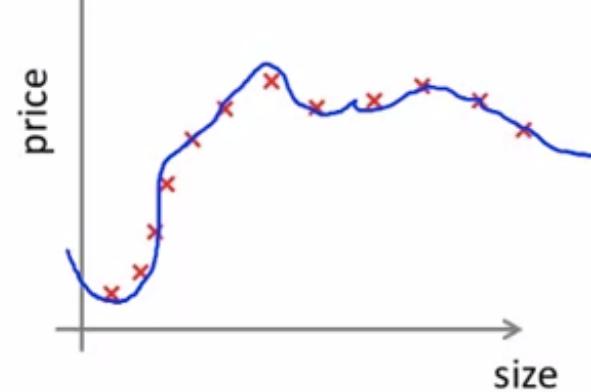
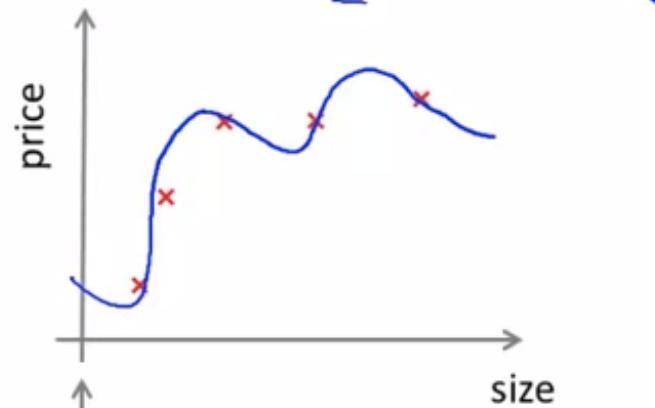


Model too complex

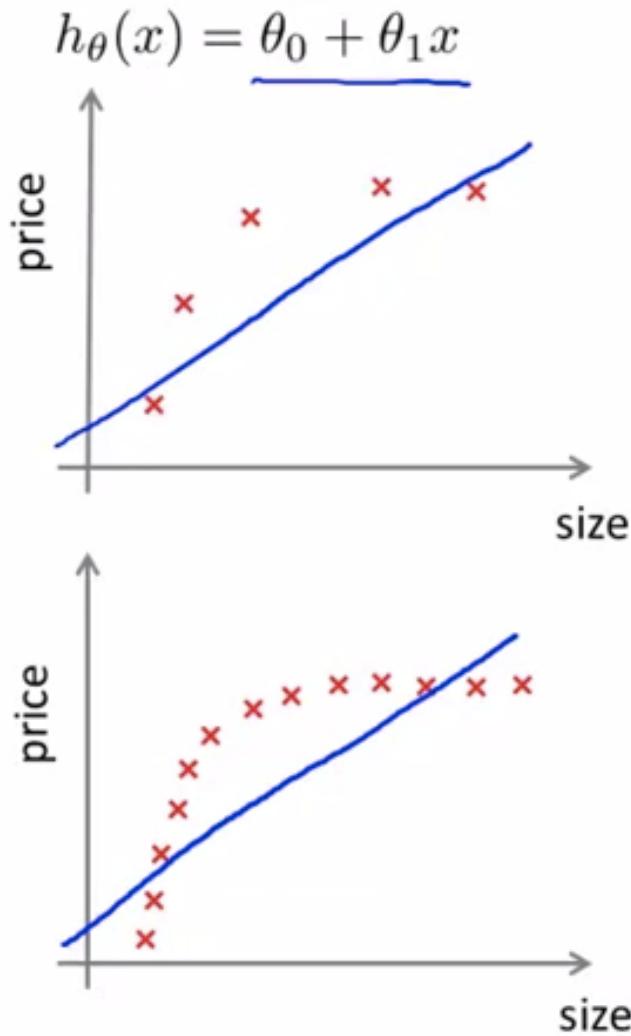
What's
the
problem?

$$h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_{100} x^{100}$$

(and small λ)



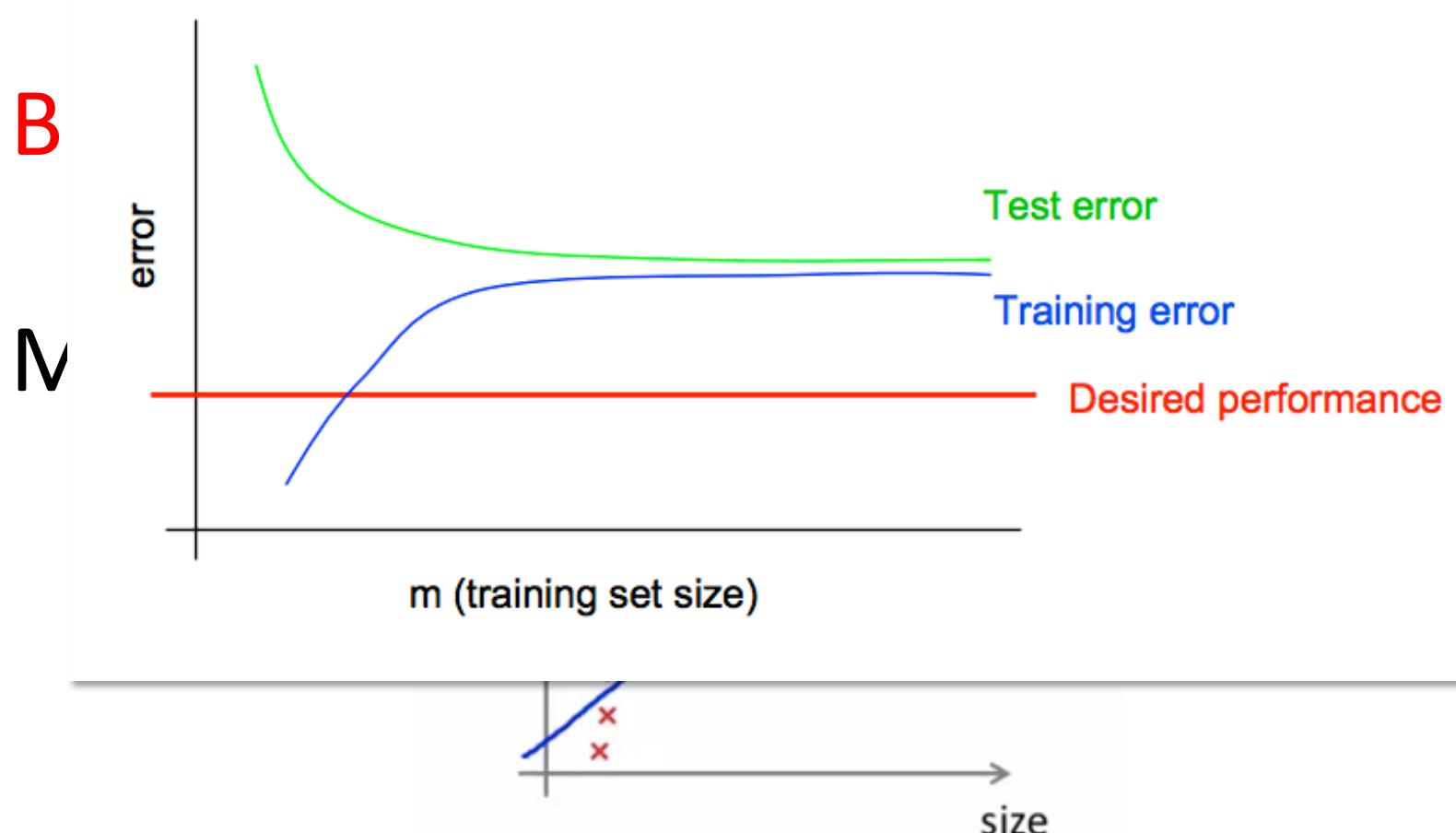
What's the problem? – Bias



Adapted from Andrew Ng – <https://www.youtube.com/watch?v=DYCv5e0lsow>, <http://see.stanford.edu/materials/aimlcs229/ML-advice.pdf>

Learning curve – Bias

$$h_{\theta}(x) = \underline{\theta_0 + \theta_1 x}$$

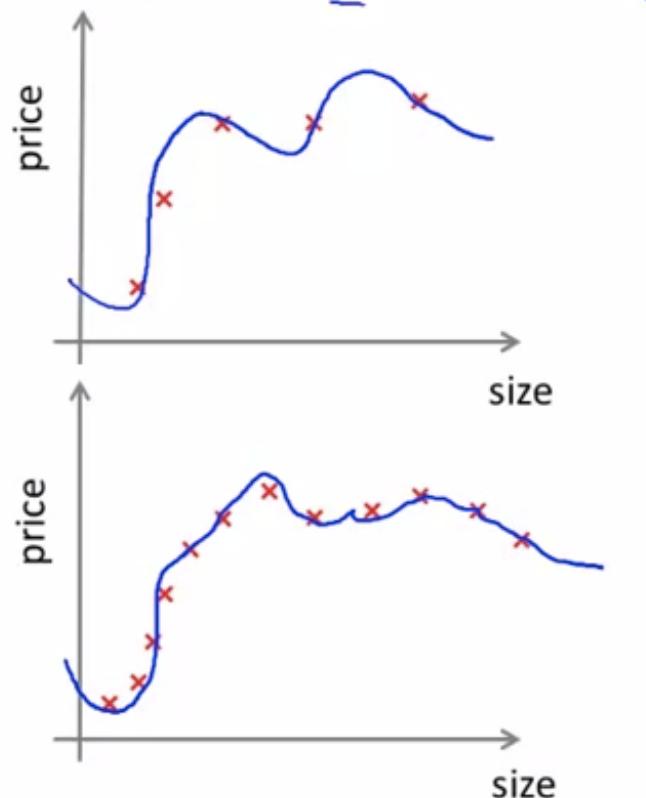


Adapted from Andrew Ng – <https://www.youtube.com/watch?v=DYCv5e0Isow>, <http://see.stanford.edu/materials/aimlcs229/ML-advice.pdf>

What's the problem? – Variance

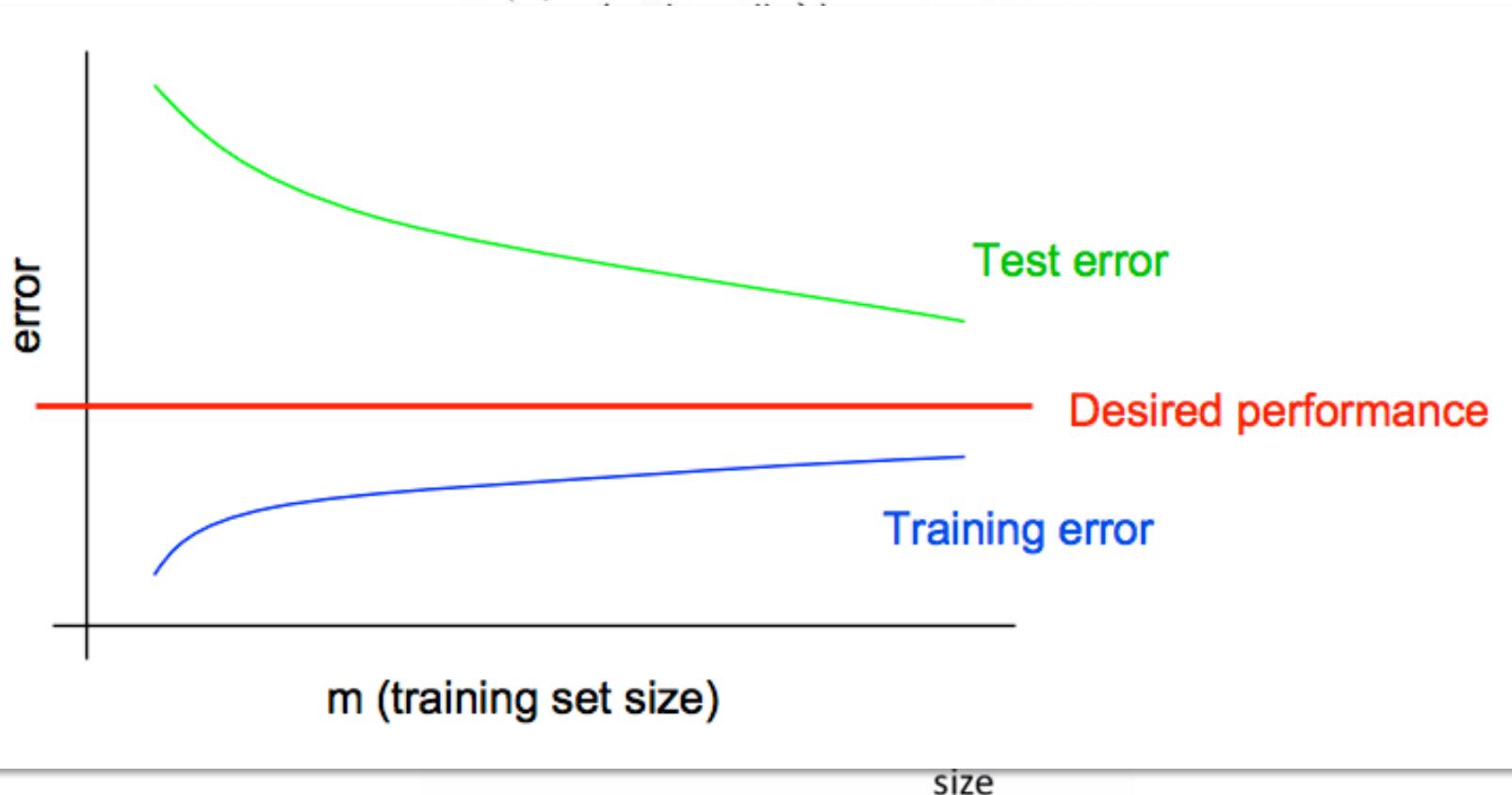
$$h_{\theta}(x) = \theta_0 + \theta_1 x + \cdots + \theta_{100} x^{100}$$

(and small λ)



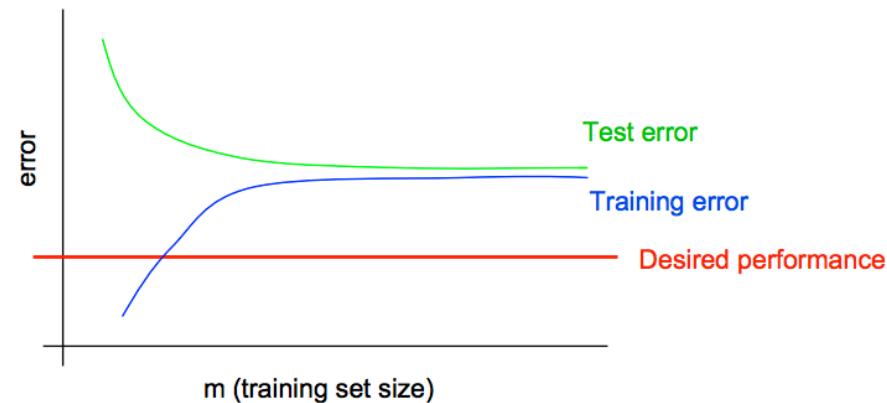
Learning curve – Variance

$$h_{\theta}(x) = \theta_0 + \theta_1 x + \cdots + \theta_{100} x^{100}$$

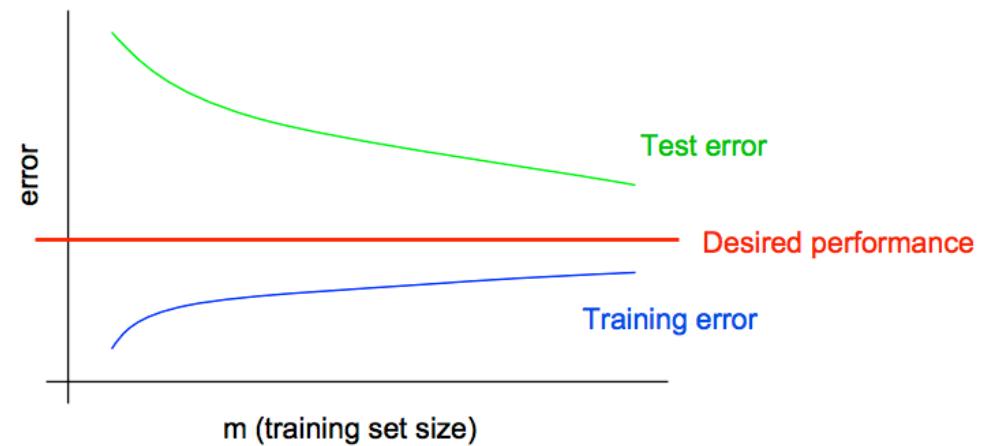


What is the next step?

Bias

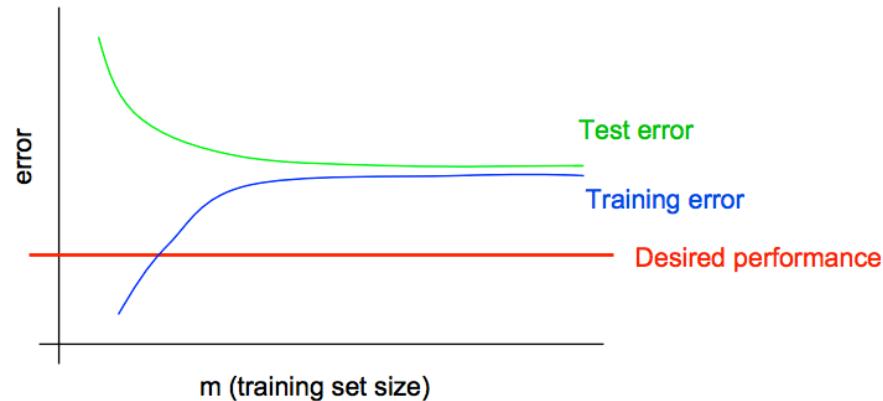


Variance



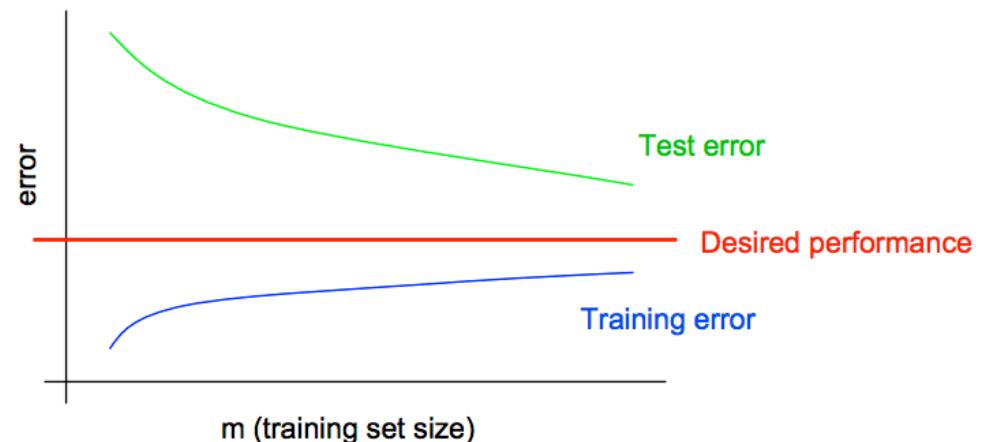
What is the next step?

Bias



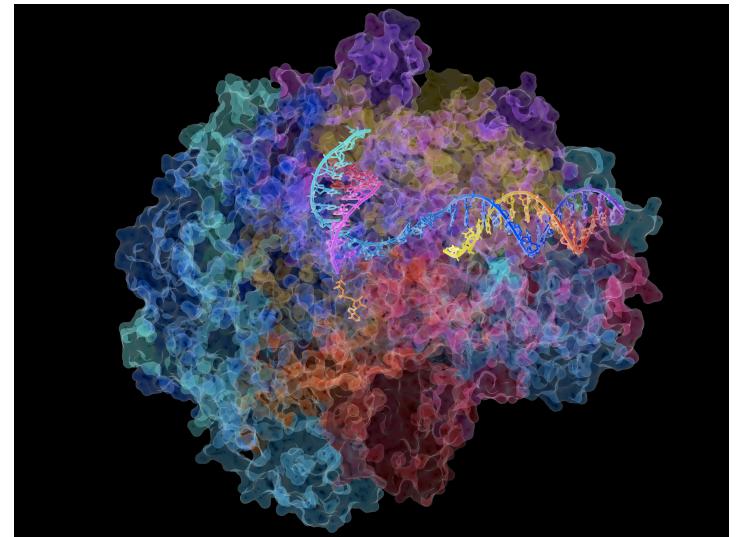
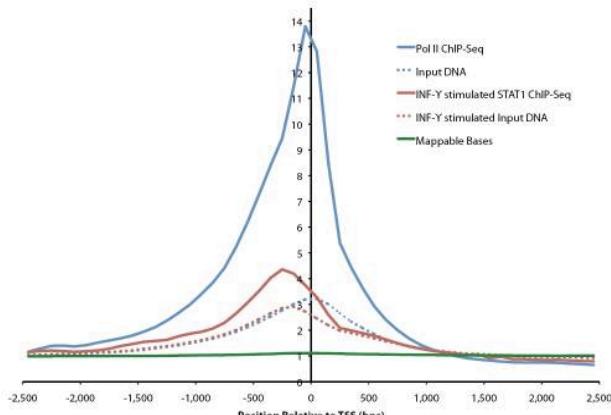
More training **features**
Train more complicated model

Variance

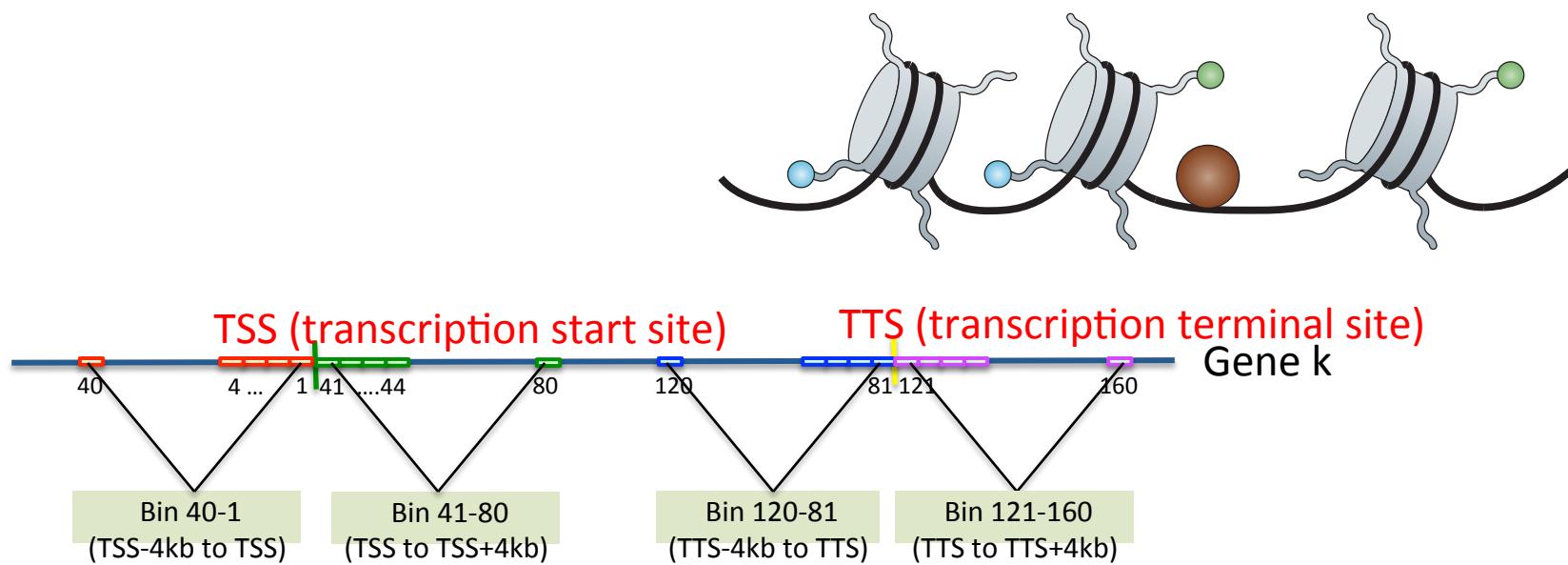


More training **examples**
Try fewer features → Dimension Reduction
Simplify model

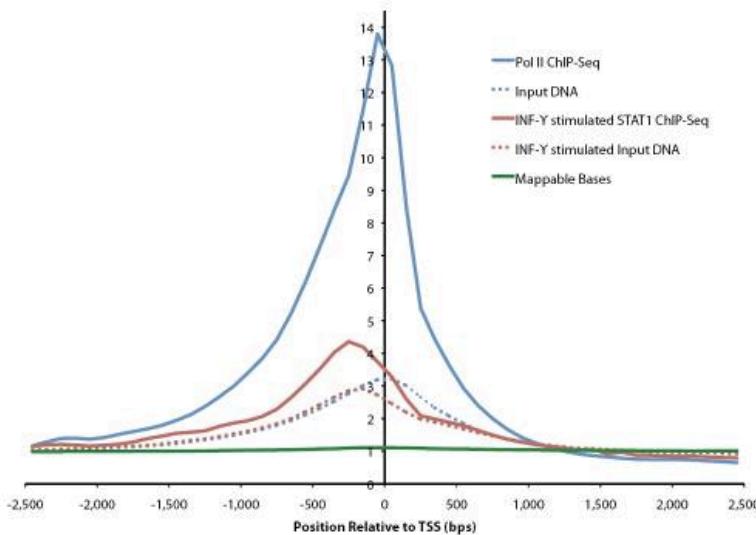
Practical Application: Predicting gene expression from ChIP-Seq signals



Where should we start?

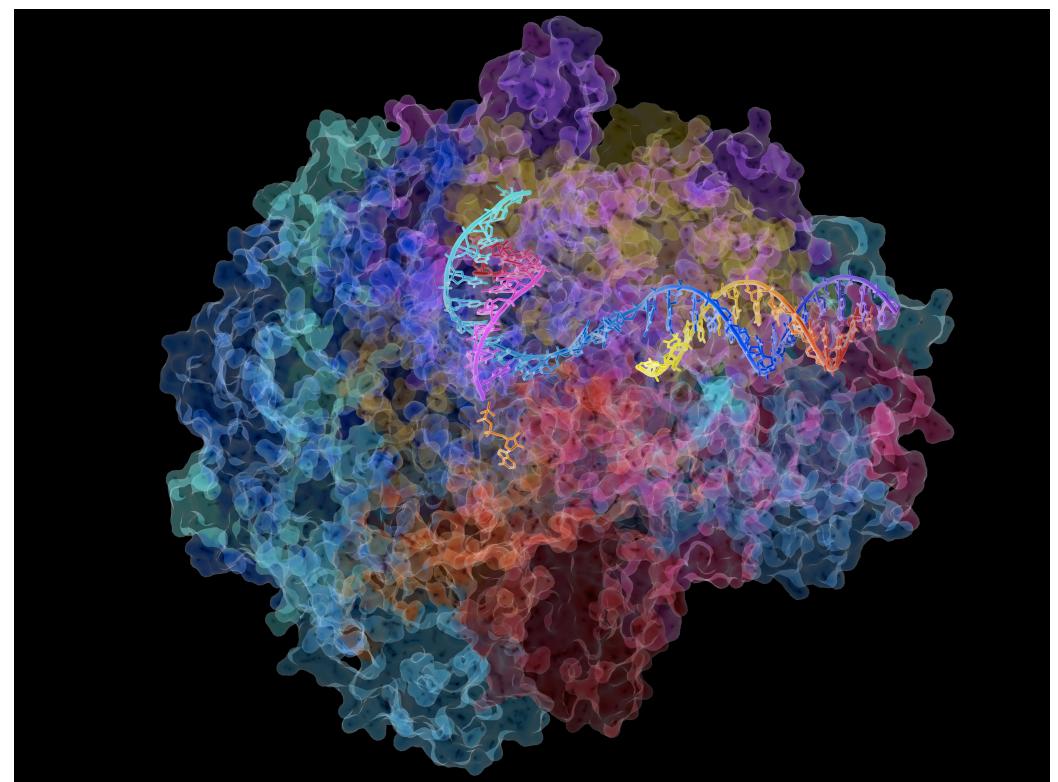


RNA is transcribed by RNA polymerase



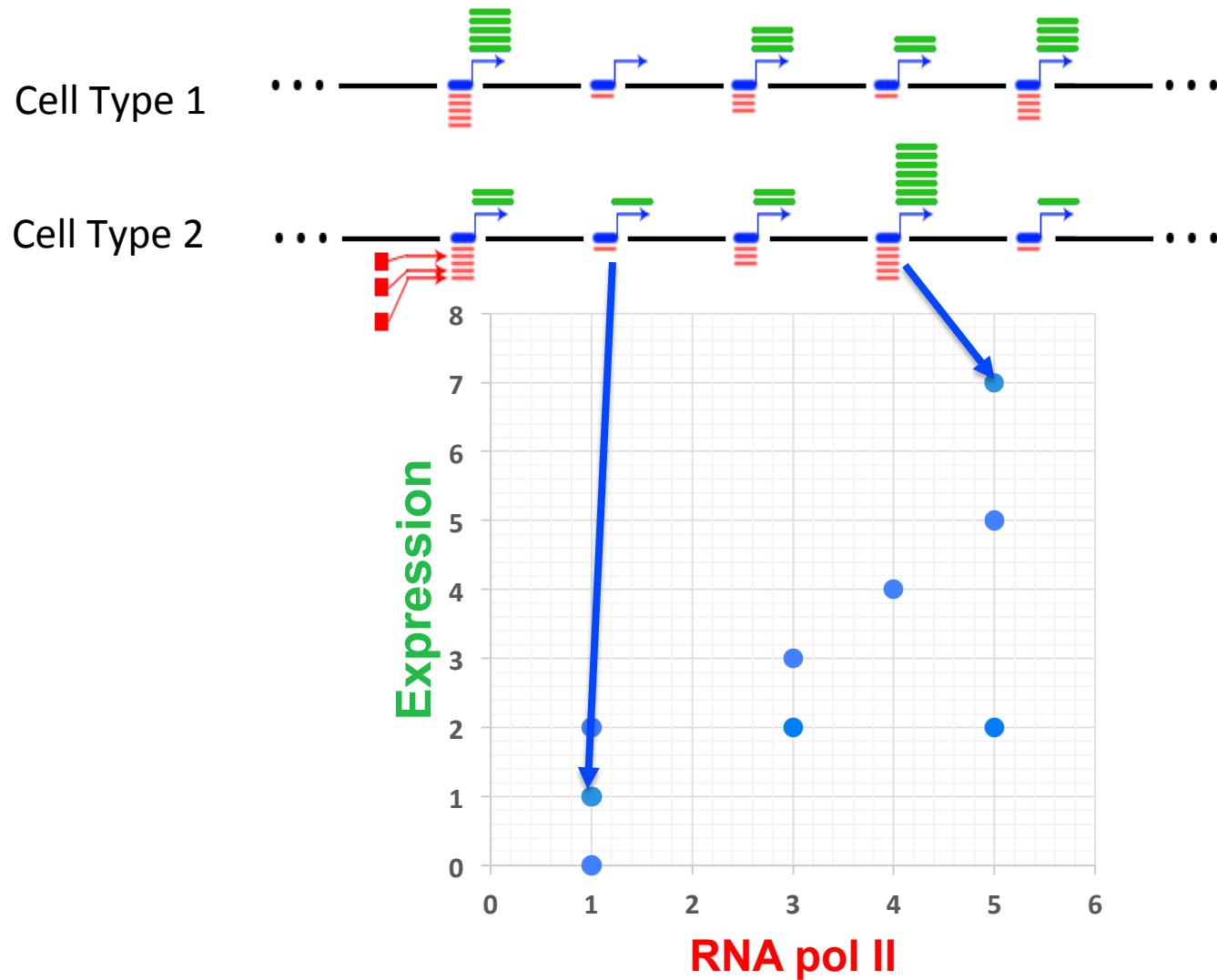
RNA pol II ChIP at
Transcription Start Sites

Rozowsky *et al.* *Nature Biotech* 2009

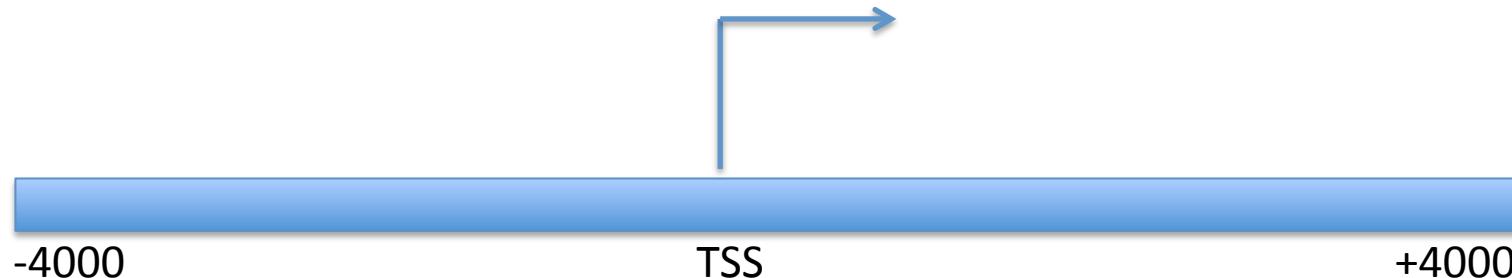


RNA polymerase II – Crystal structure
Roger Kornberg Nobel Prize

Relating Genomic Inputs to Outputs



Initial model



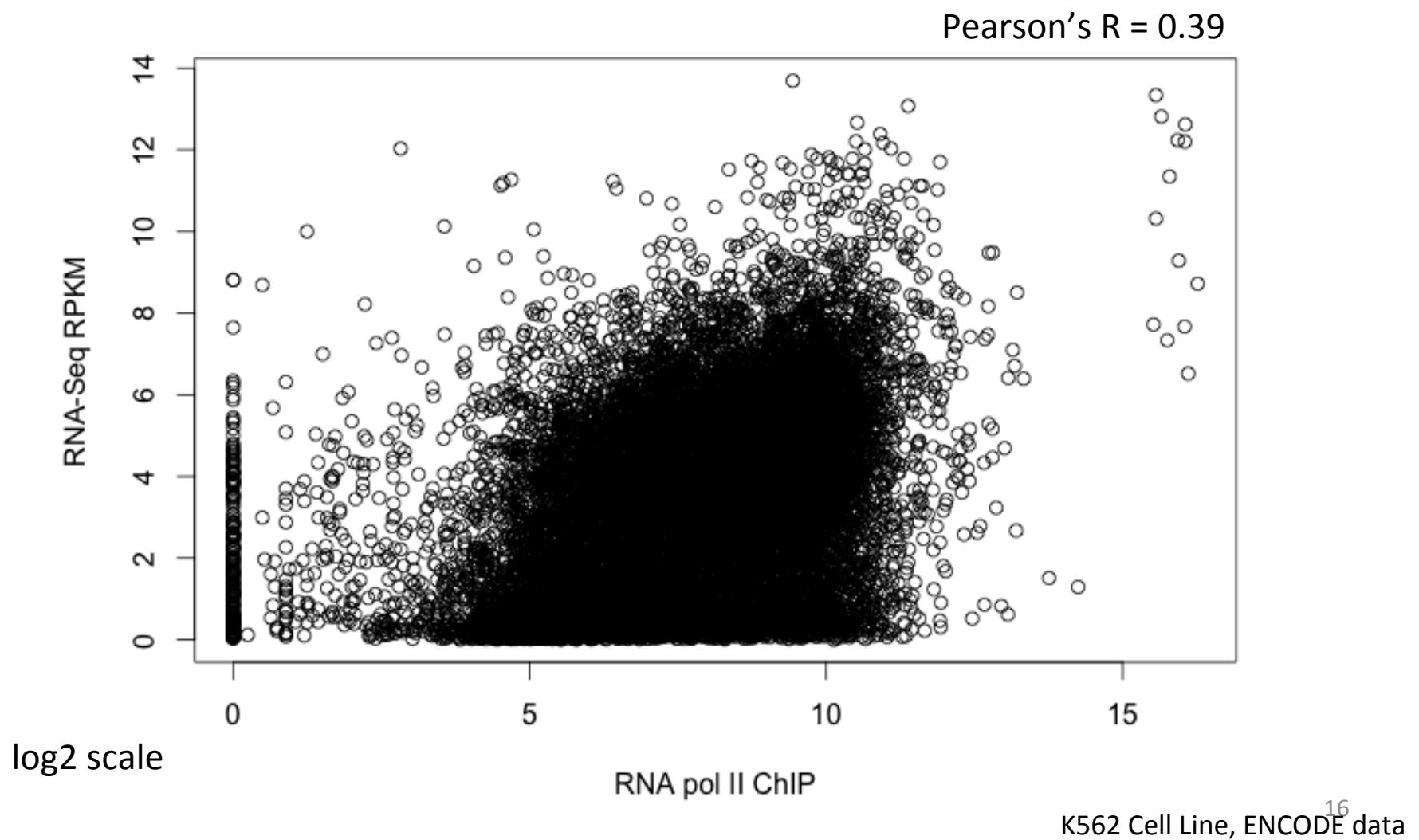
data from K562 cell line, ENCODE consortium

Sum pol II ChIP signal across 8000 bp centered around transcription start site.

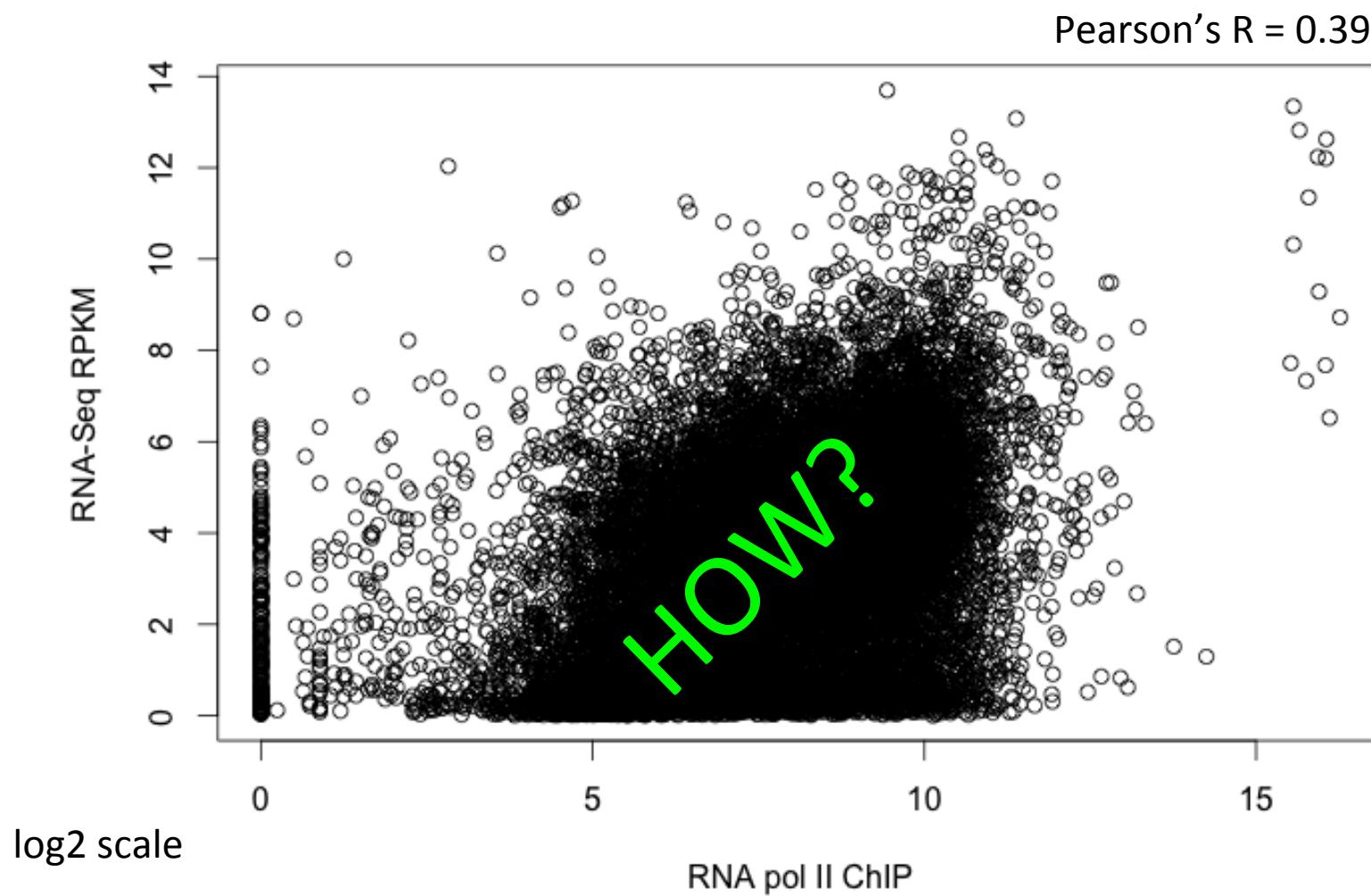
$$\log_2(\text{RNA-Seq RPKM}) = a + b * \log_2(\text{RNA Pol II ChIP})$$

Why log scale??

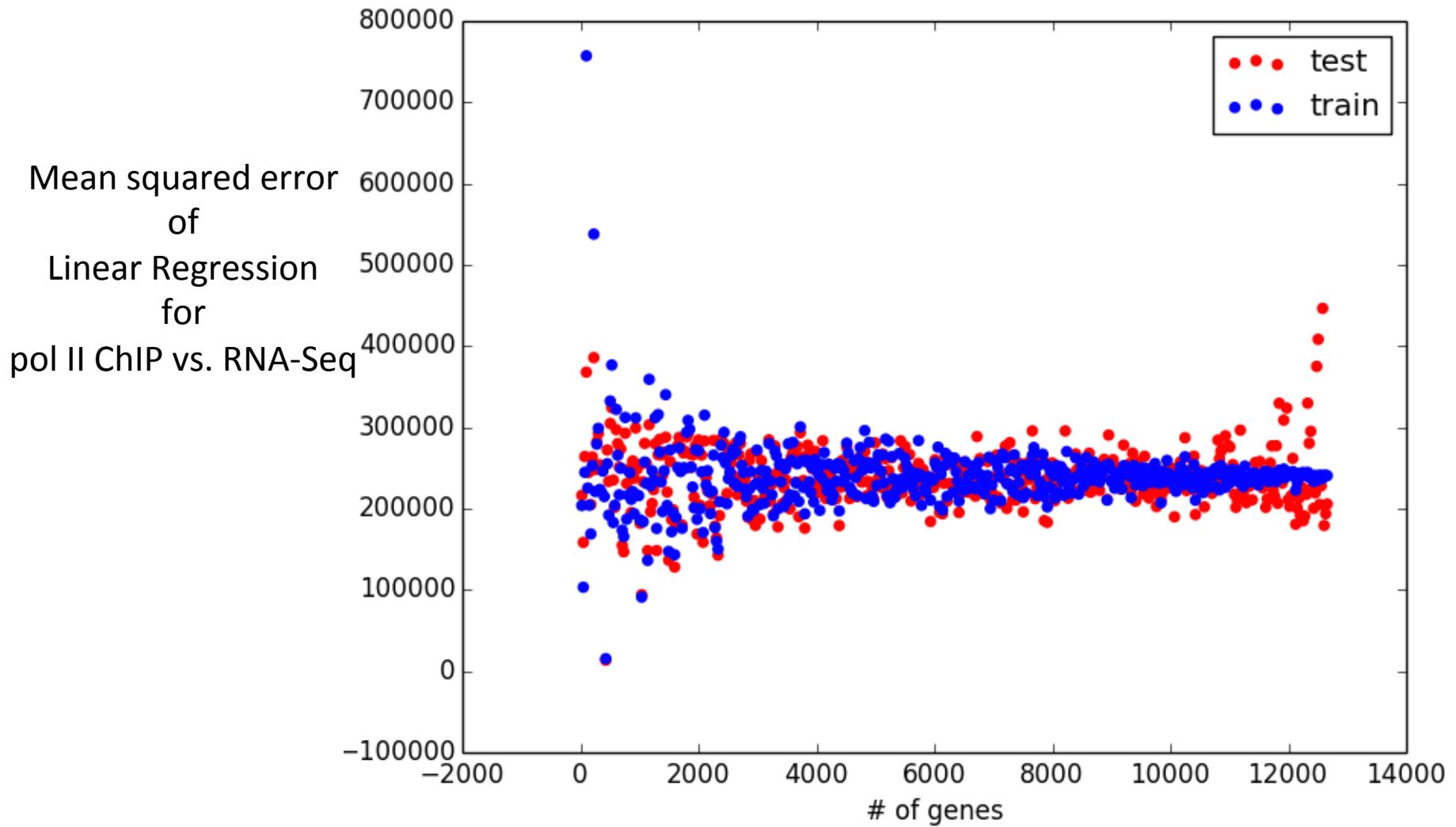
Can we do better than this?



Can we do better than this?



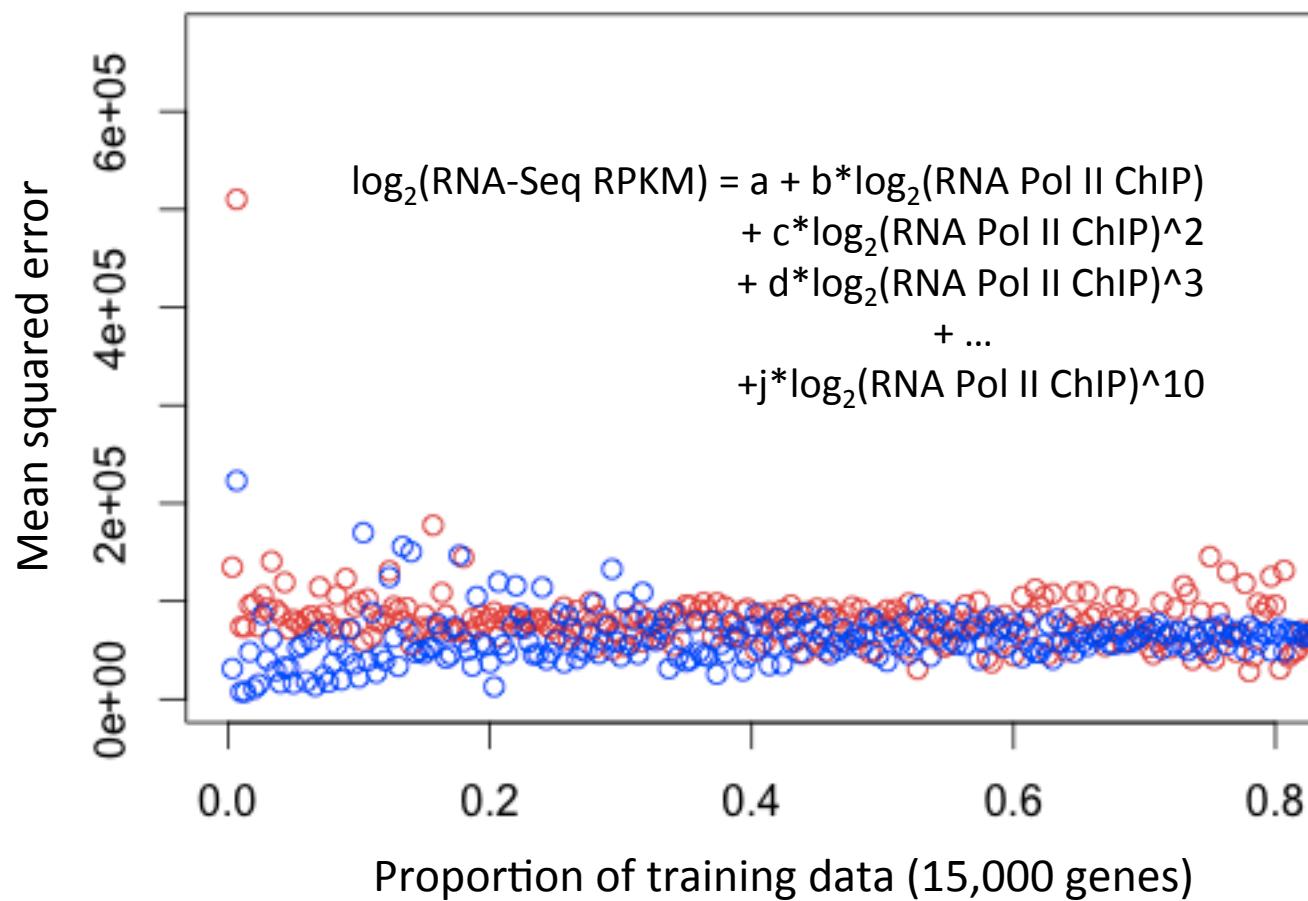
Learning curve: What's our problem?



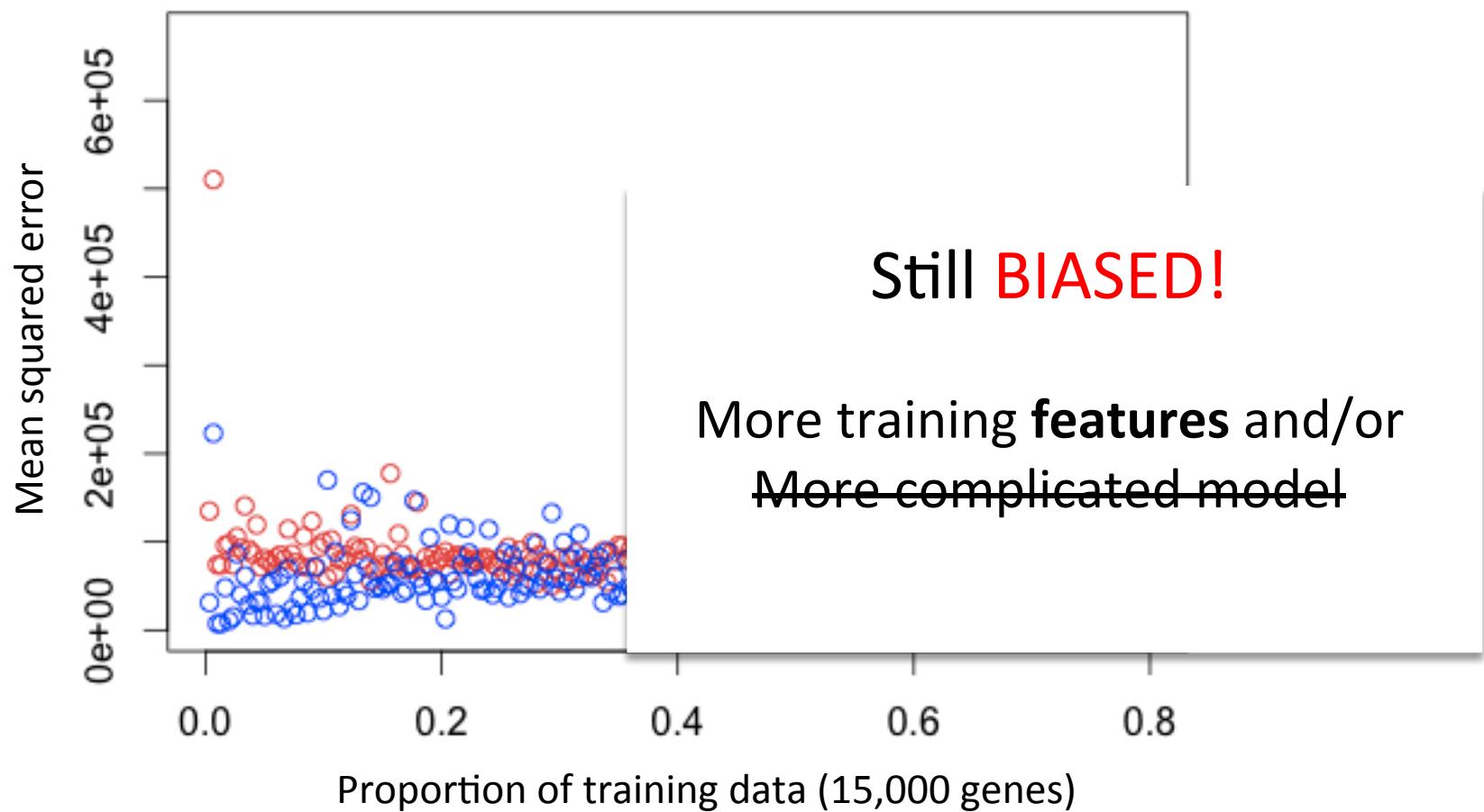
Learning curve: What's our problem?

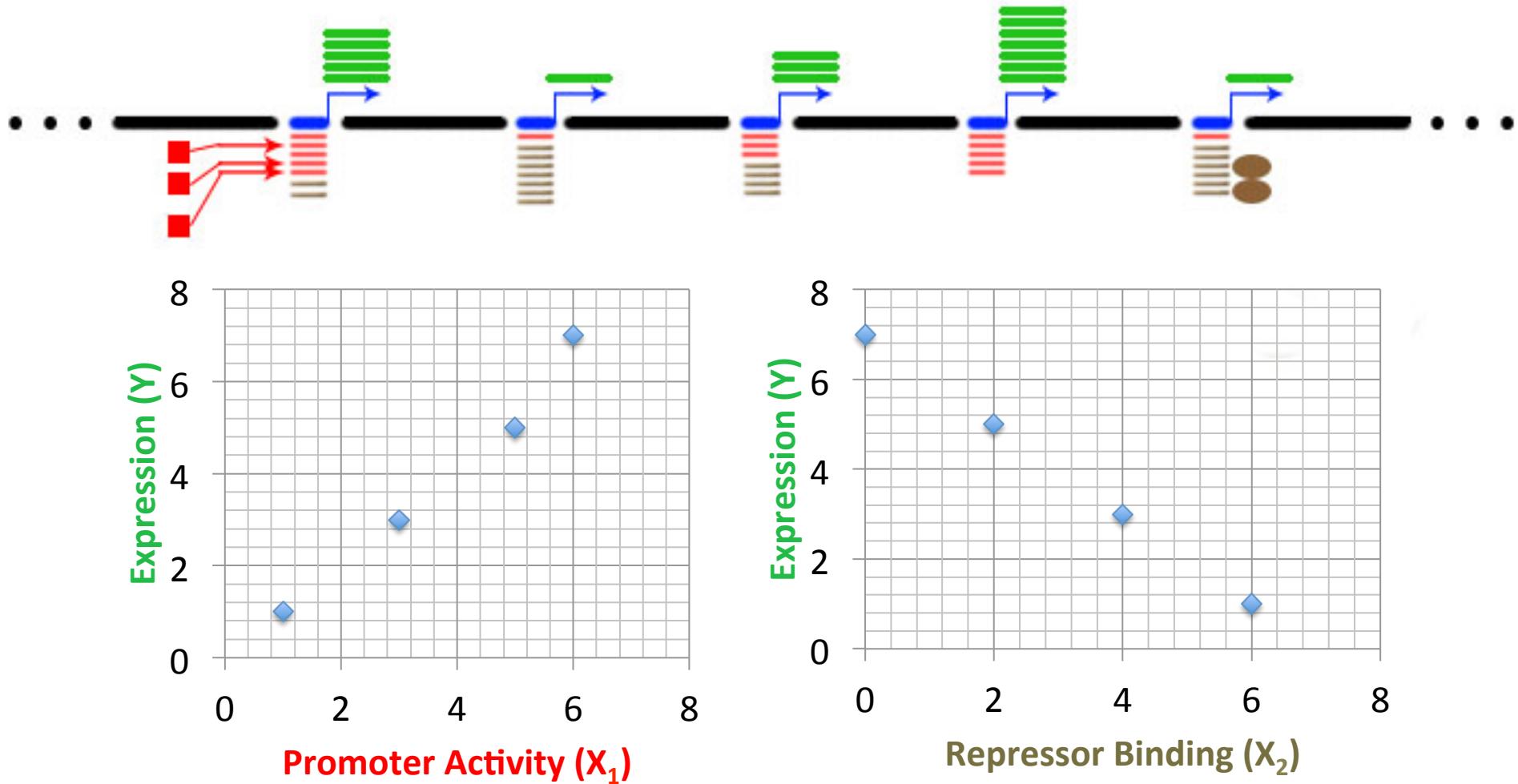


Polynomial regression: polII ChIP vs. RNA-Seq



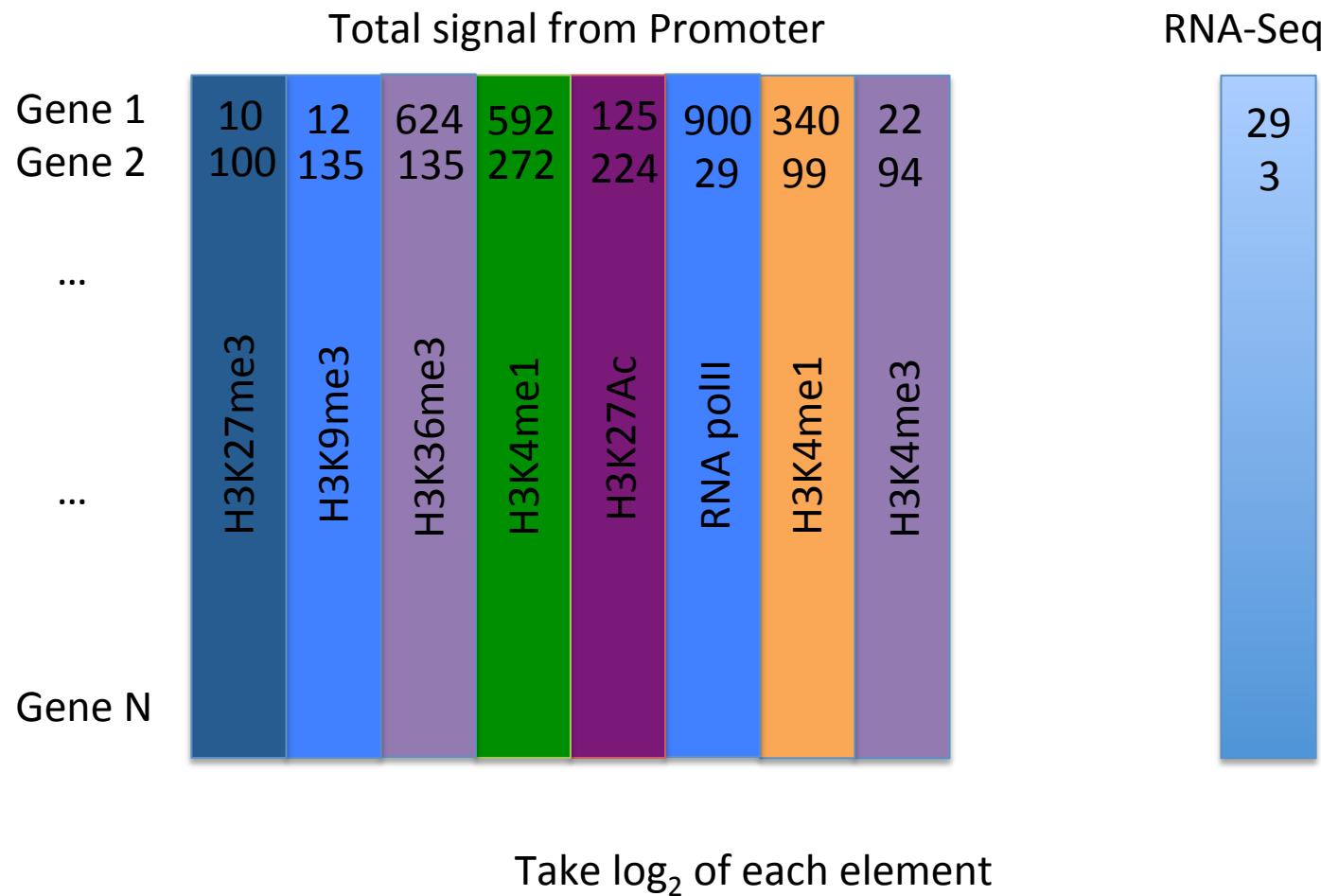
Polynomial regression: polII ChIP vs. RNA-Seq



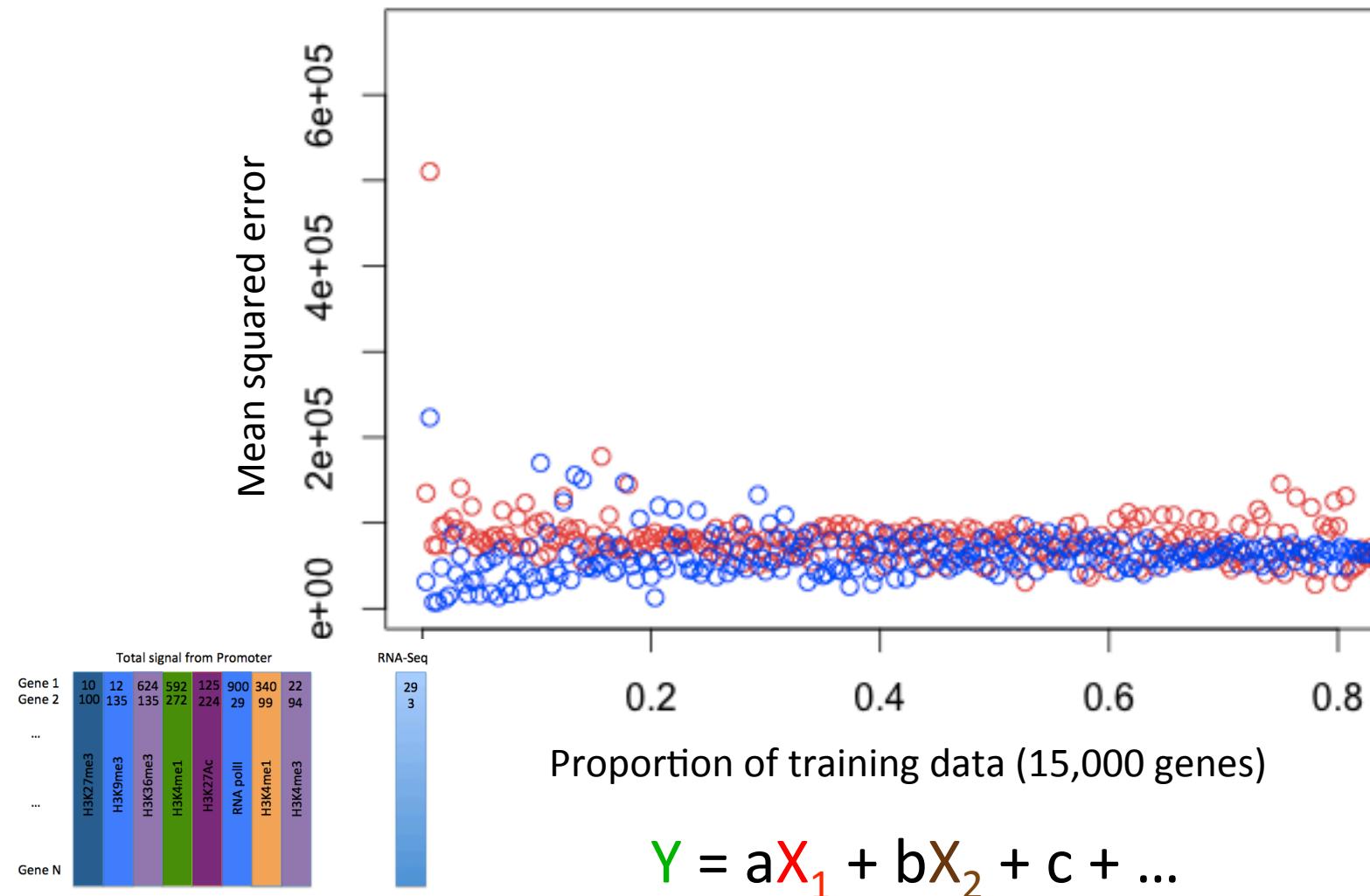


$$Y = aX_1 + bX_2 + c$$

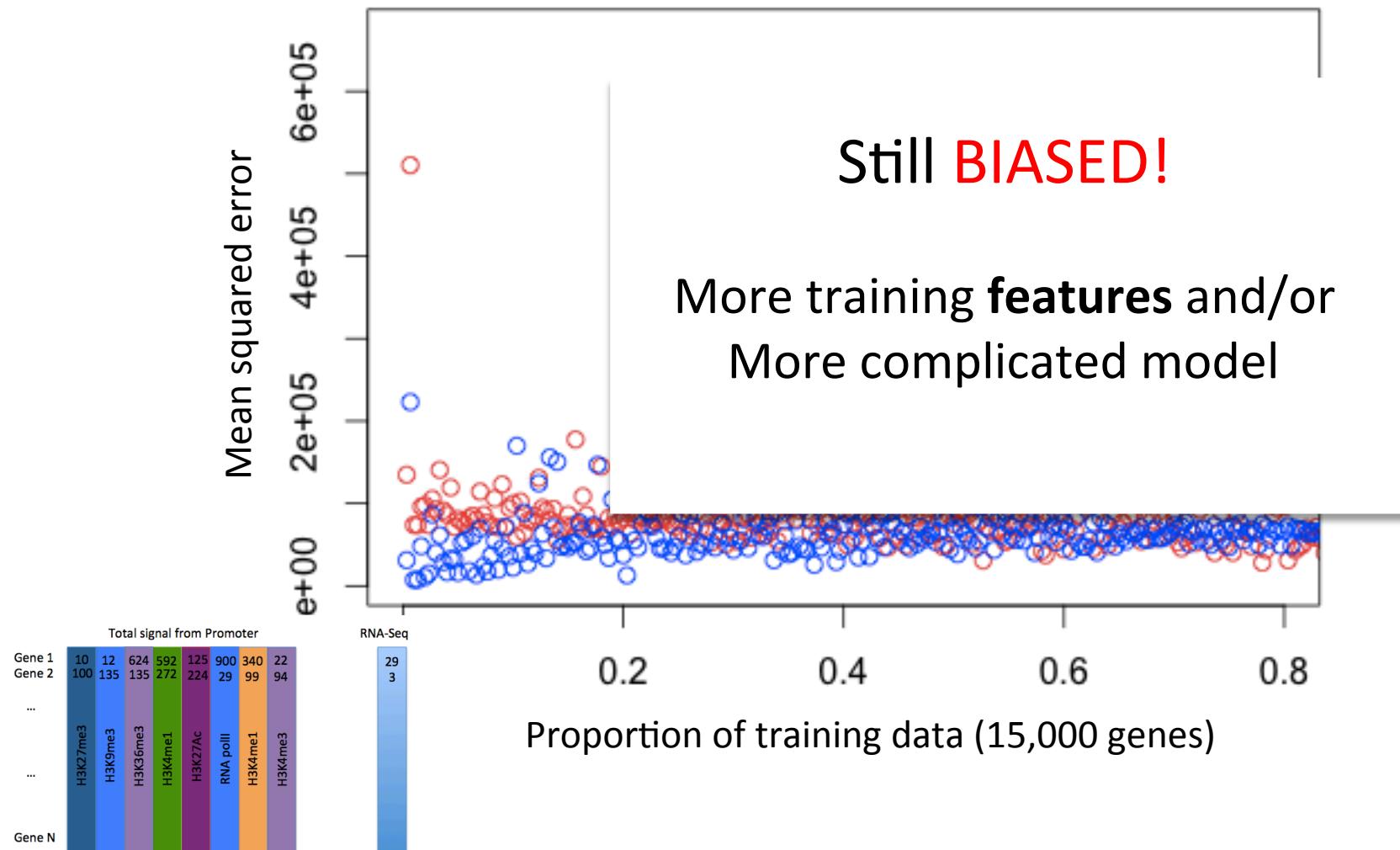
Adding more signals



Multiple Linear Regression



Multiple Linear Regression



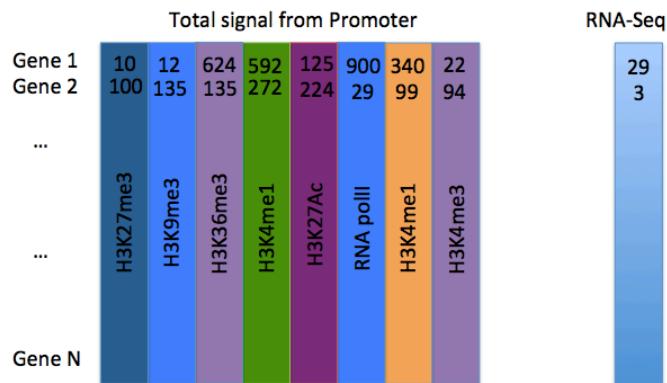
Random Forest Regression

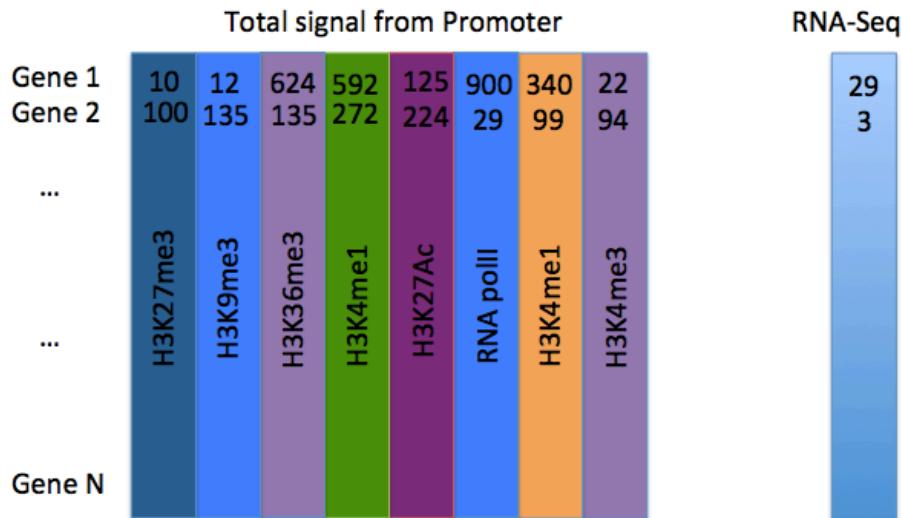
Log-transformed

Training correlation: 0.93
Test correlation: 0.52

Not Log-transformed

Training correlation: 0.95
Test correlation: 0.16

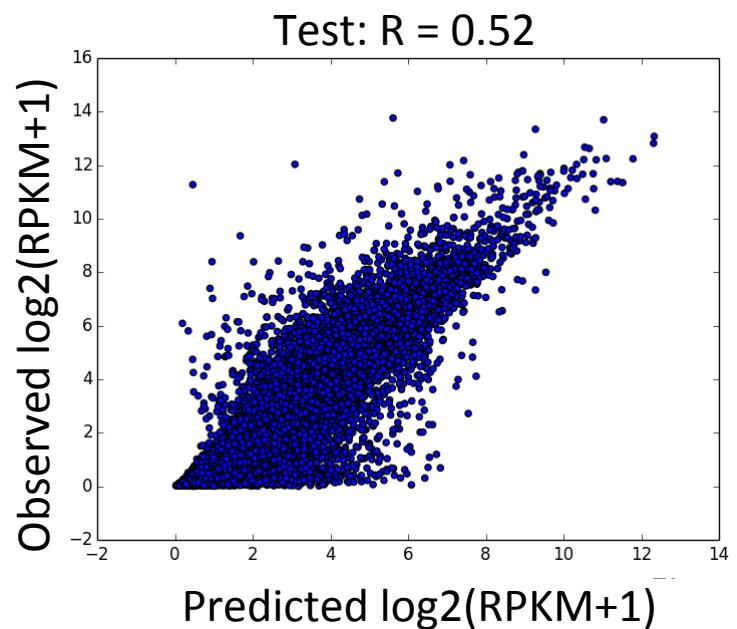
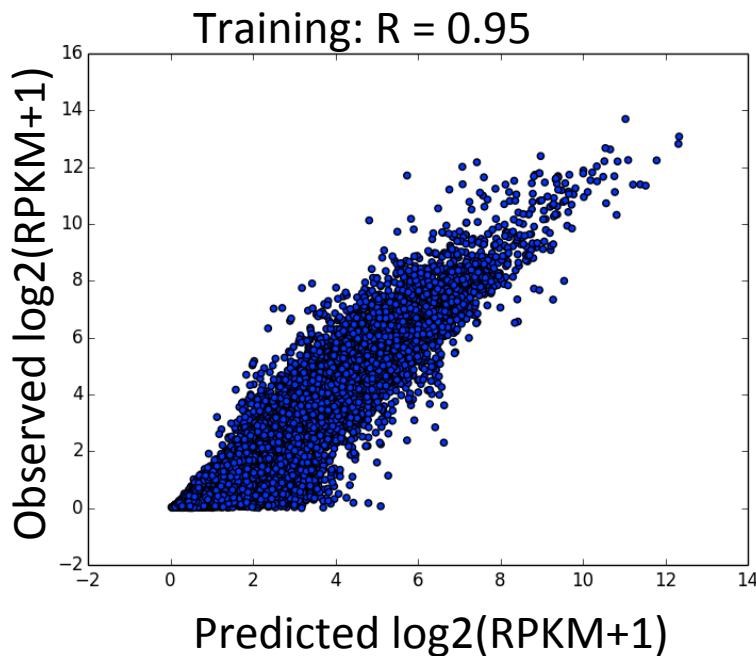




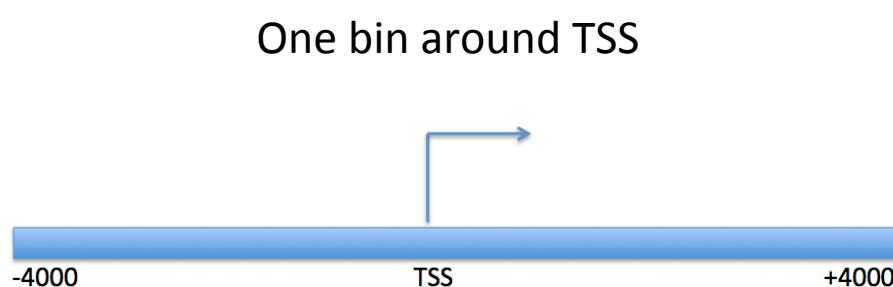
Random Forest Regression

Candy:

What's the problem – Bias or Variance?
What should we do now?

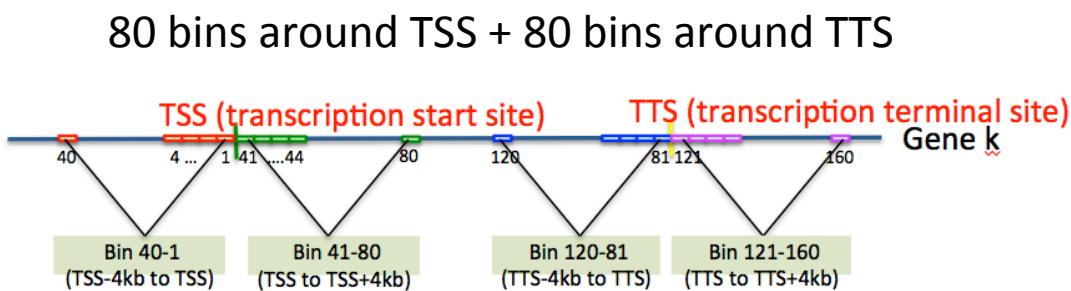


What's the best model setup?

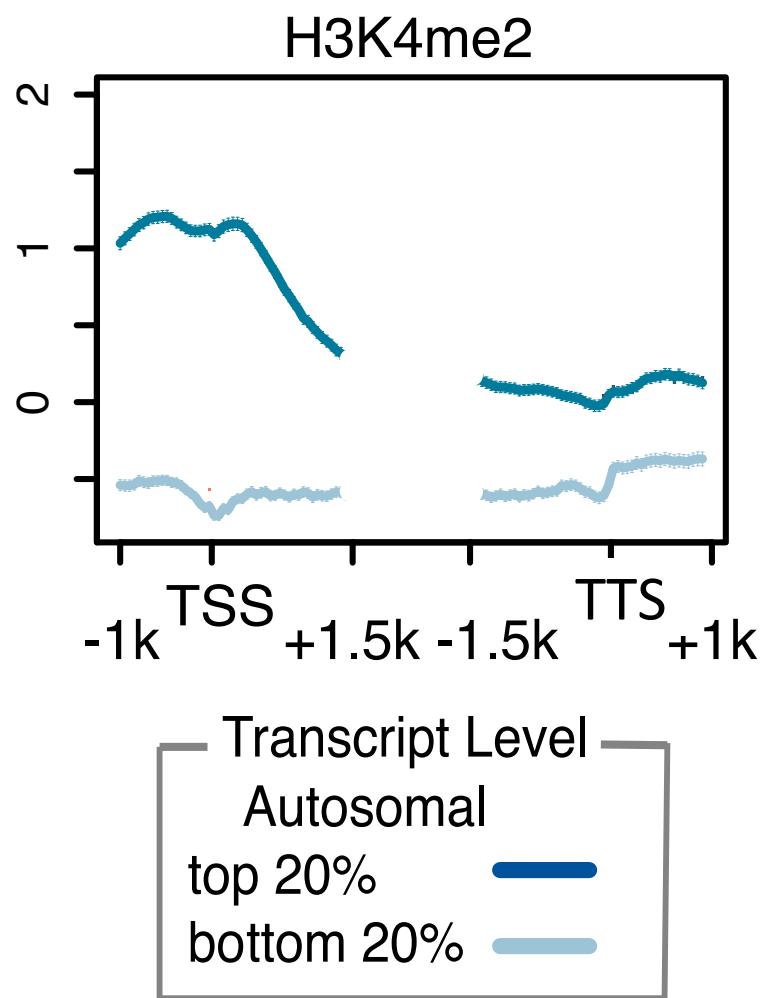


Candy: Which setup do we expect to perform better?

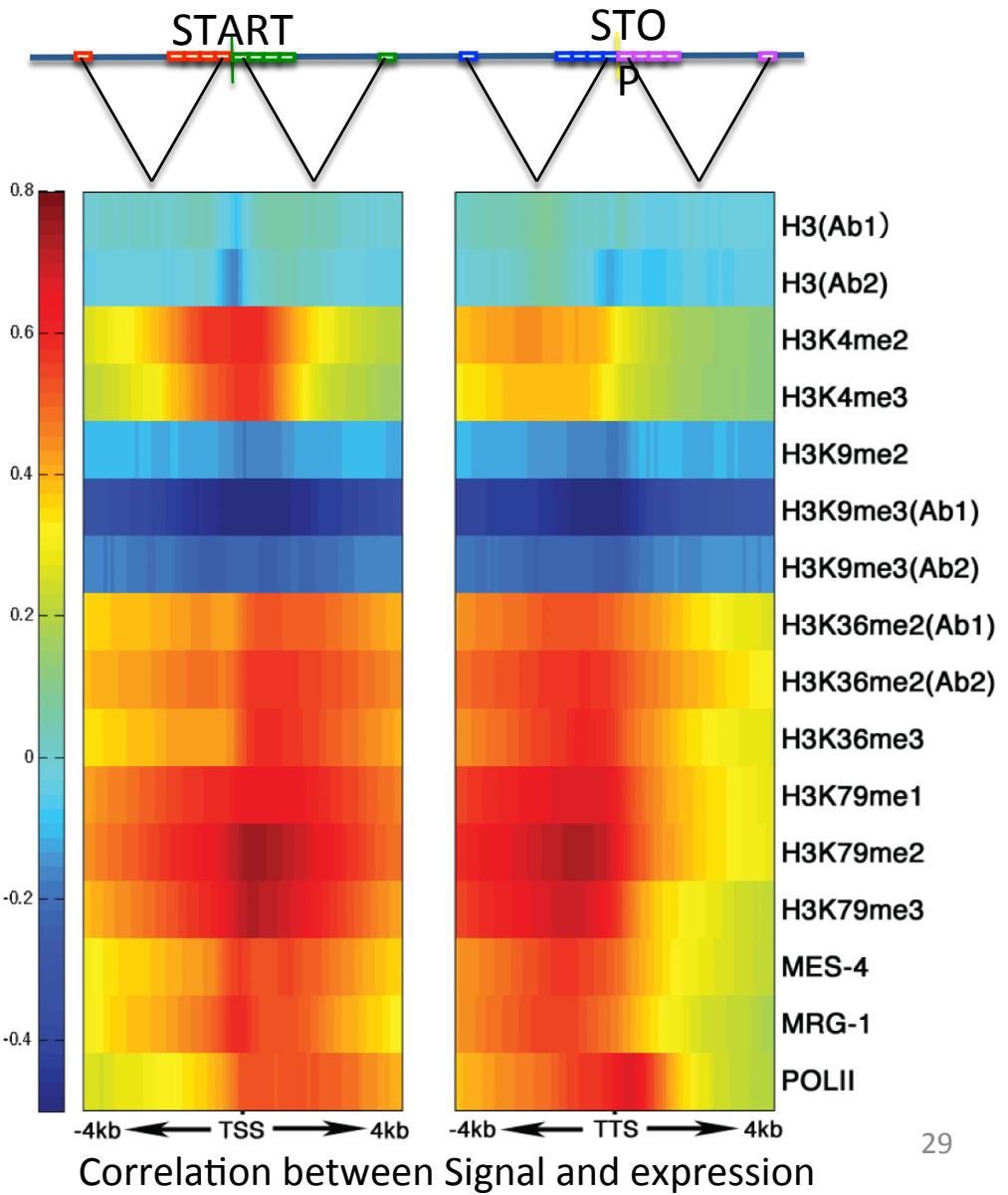
Vs.



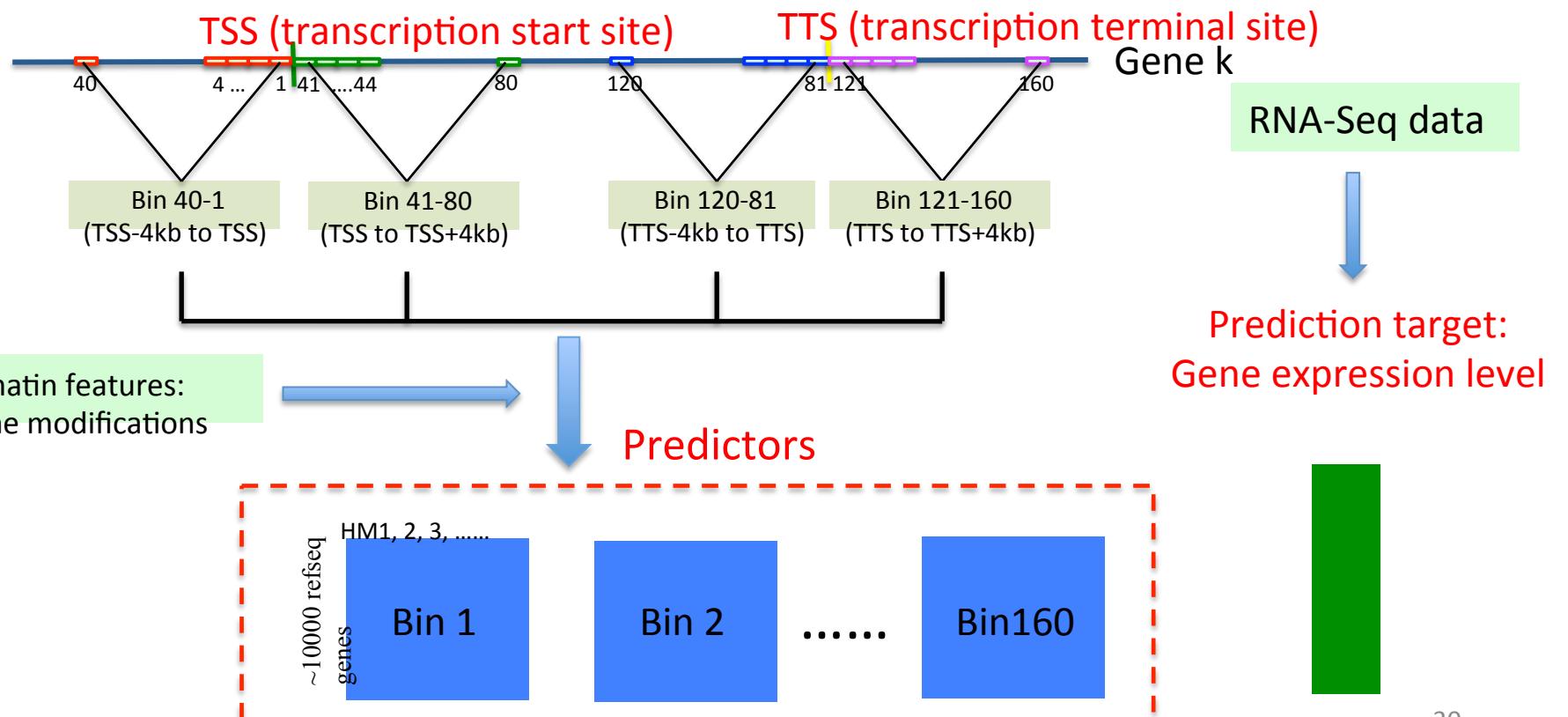
Effects of signal depend on Location!



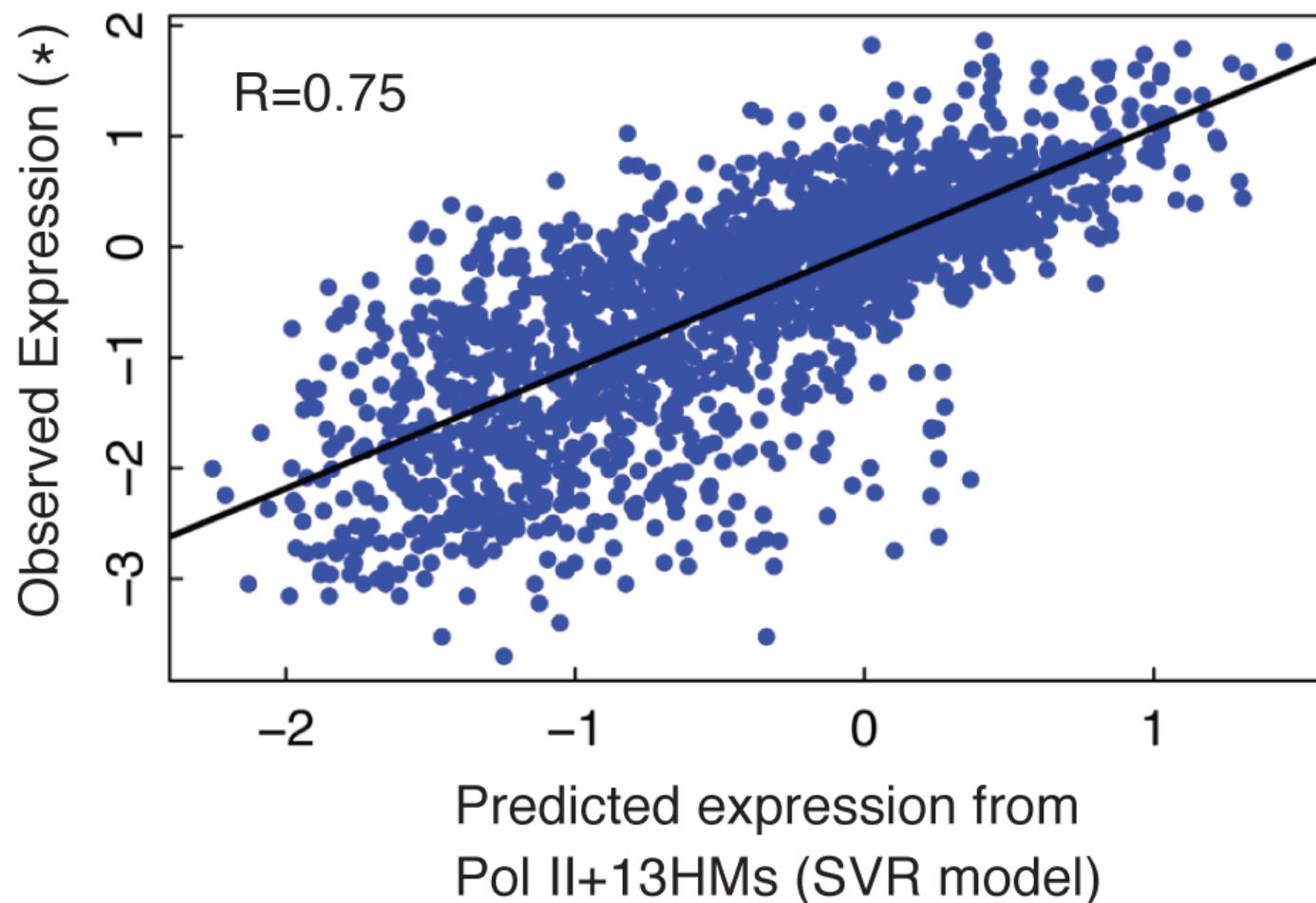
Gerstein*,..., Cheng* et al. 2010,
Science



Setting up the model



Support vector regression to predict gene expression levels

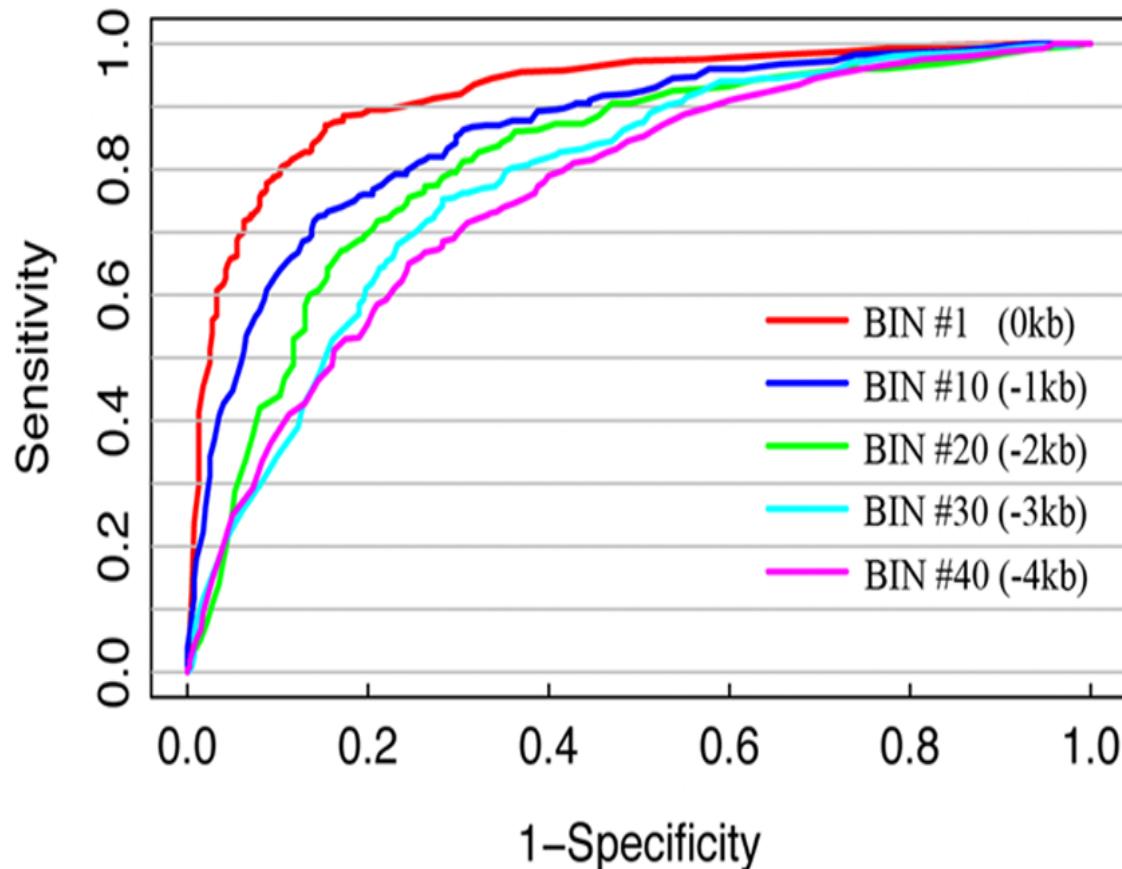


Context (TA bias)

- My implementation:
 - Train correlation = 0.95
 - Test correlation = 0.58

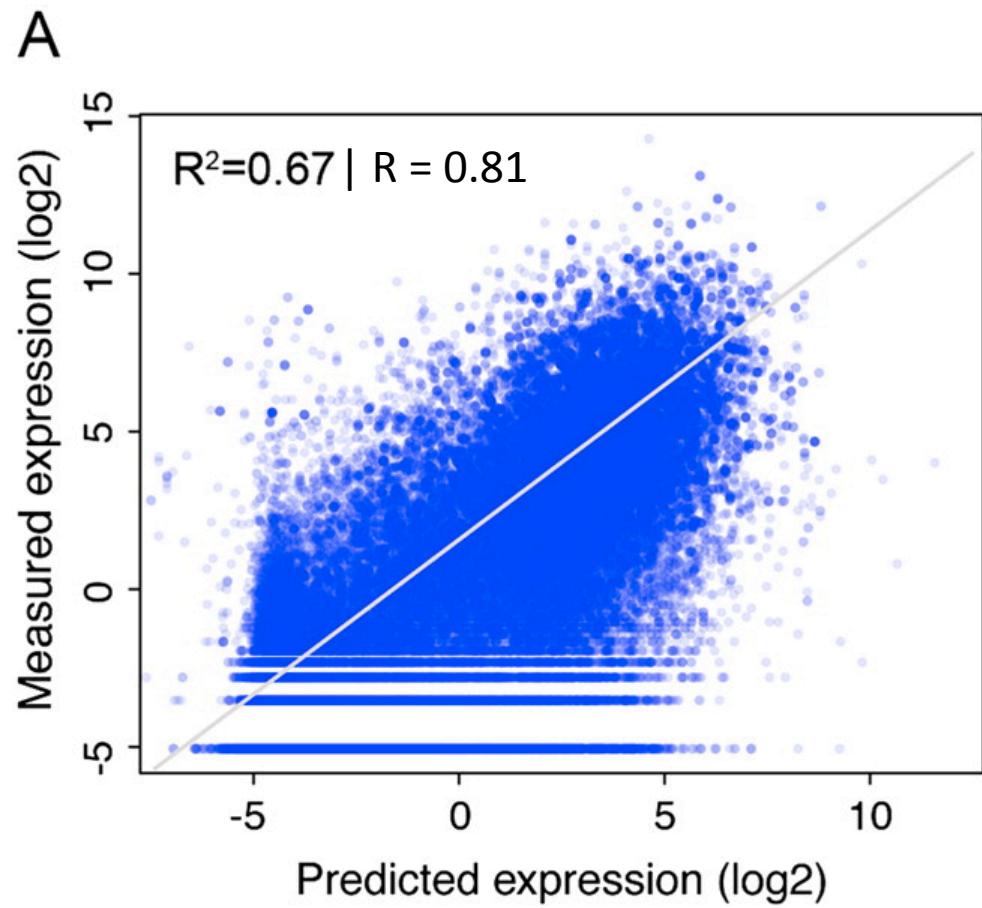
Support vector machine to classify genes with high, medium and low expression

Candy:
Describe
a ROC
curve!

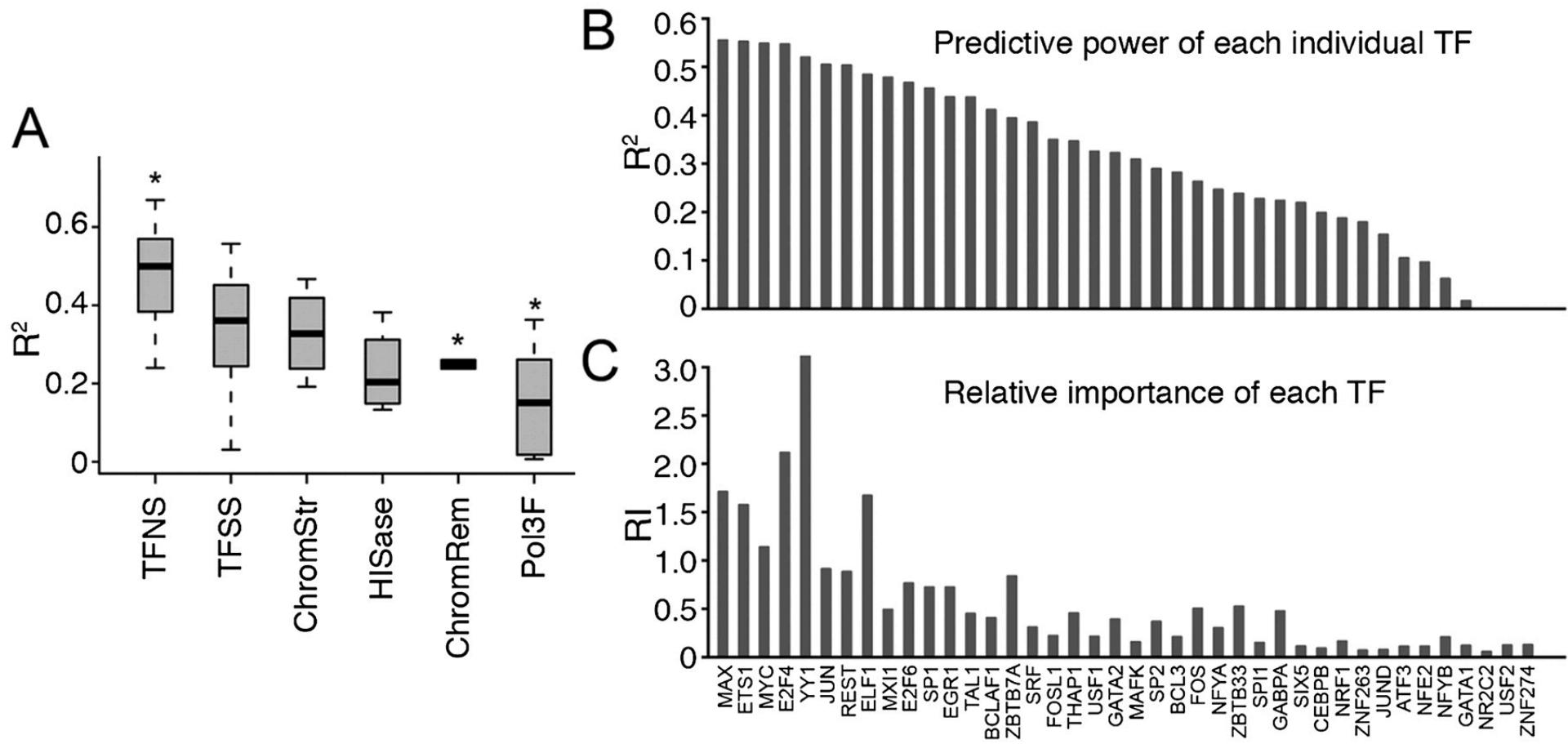


- ✓ Areas close to TSS predict expression better

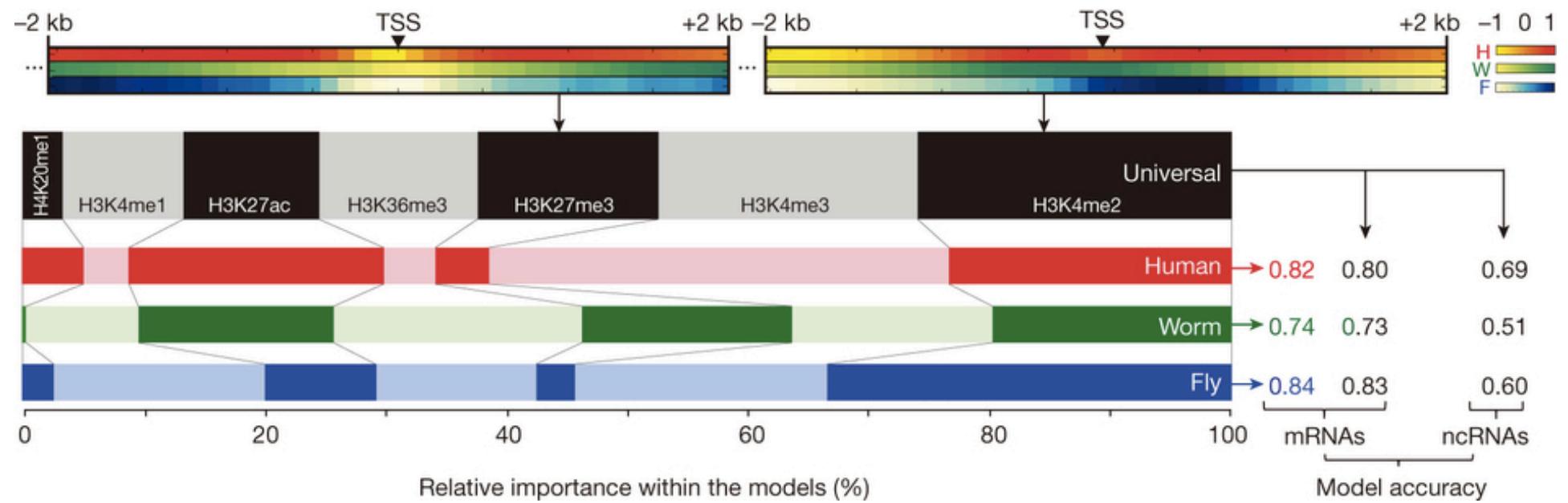
Predicting Gene Expression with Transcription Factor ChIP-Seq signals



Predicting Gene Expression with Transcription Factor ChIP-Seq signals



Modeling Transcription Between Organisms



Why do we care?

- What are the benefits of a quantitative model?
- Does this model help us understand the mechanism of transcription?

For discussion

- Will the prediction model perform accurately in cells with a transcription factor knocked out?

