# Advanced Machine Learning

## Applications:
## Computational Biology

**Eric Xing**

**Lecture 19, August 14, 2008**

**Reading: see class homepage**

---

# Biological Data Analysis

- Dynamic, noisy, heterogeneous, high-dimensional data

- "High-resolution" inference
- Parsimonious
- Scalability
- Stability
- Sample complexity
- Confidence bound

## Structured Prediction Problem

- Unstructured prediction

$$\mathbf{x} = \begin{pmatrix} x_{11} & x_{12} & \cdots \\ x_{21} & x_{22} & \cdots \\ \vdots & \vdots & \cdots \end{pmatrix} \qquad \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \end{pmatrix}$$

- Structured prediction
  - Part of speech tagging

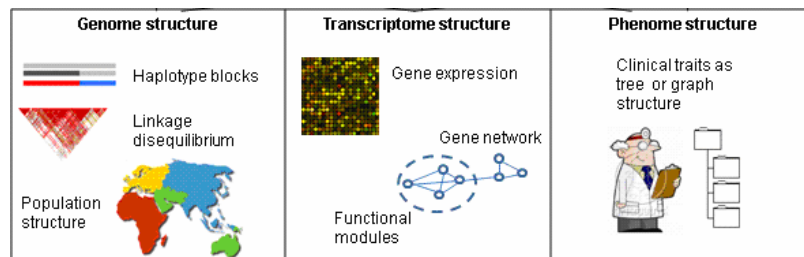    $\mathbf{x} =$ "Do you want sugar in it?" $\Rightarrow$ $\mathbf{y} =$ <verb pron verb noun prep pron>

  - Image segmentation

    $$\mathbf{x} = \begin{pmatrix} x_{11} & x_{12} & \cdots \\ x_{21} & x_{22} & \cdots \\ \vdots & \vdots & \cdots \end{pmatrix} \qquad \mathbf{y} = \begin{pmatrix} y_{11} & y_{12} & \cdots \\ y_{21} & y_{22} & \cdots \\ \vdots & \vdots & \cdots \end{pmatrix}$$

## Genome and Phenome Structures



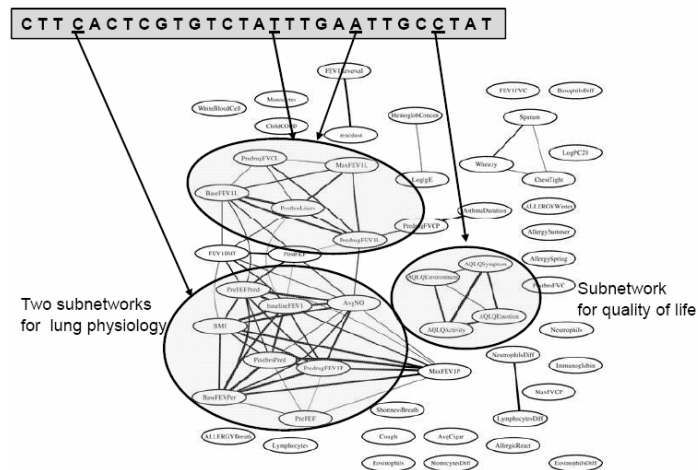| Genome structure | Transcriptome structure | Phenome structure |
|---|---|---|
| Haplotype blocks | Gene expression | Clinical traits as tree or graph structure |
| Linkage disequilibrium | Gene network | |
| Population structure | Functional modules | |

## Genome-Phenome Association

| Traditional view | Modern view |
|---|---|
| **causal SNP** | **causal SNP networks** |
| **ACGTTTTACTGTACAATT** | **ACGTTTTACTGTACAATT** |
| **a univariate phenotype:** | **Multivariate complex syndrome (e.g., asthma):** |
| **i.e., disease/control** | **age at onset,** |
| | **presence/absence of eosinophilic inflammation,** |
| | **history of eczema** |
| | **genome-wide expression profile** |
| | **...** |

---

## The Asthma Phenotype Network



C T T C A C T C G T G T C T A T T T G A A T T G C C T A T

Two subnetworks for lung physiology

Subnetwork for quality of life

# The Asthma Phenotype Network



Pleotropic effects

Two subnetworks for lung physiology

Subnetwork for quality of life

Epistatic effects

CTTCACTCGTGTCTATTTGAATTGCCTAT

---

# Inferring Genome-Phenome Association



Traditional view

**ACGTTTTACTGTACAATT**

- **Pair-wise association tests?**
  - Ignore SNP dependencies
  - Many many FPs

- **Regression?**
  - Over-stringent on coupled SNPs

Modern view

**ACGTTTTACTGTACAATT**

- **Structured regularized regression**
  - ✓ explicitly capture structures
  - ✓ efficient
  - ✓ sparse (parsimonious)
  - ✓ provable guarantees

## Association Mapping

| | Phenotype (BMI) | Genotype |
|---|---|---|
| Individual 1 | 2.5 | . . C . . . . . T . . C . . . . . . . T . . . <br> . . C . . . . . A . . C . . . . . . . T . . . |
| Individual 2 | 4.8 | . . G . . . . . A . . G . . . . . . . A . . . <br> . . C . . . . . T . . C . . . . . . . T . . . |
| Individual N | 4.7 | . . G . . . . . T . . C . . . . . . . T . . . <br> . . G . . . . . T . . G . . . . . . . T . . . |

Benign SNPs    Causal SNP

## Association Mapping as Regression

| | Phenotype (BMI) | Genotype |
|---|---|---|
| Individual 1 | 2.5 | . . 0 . . . . . 1 . . 0 . . . . . . . 0 . . . |
| Individual 2 | 4.8 | . . 1 . . . . . 1 . . 1 . . . . . . . 1 . . . |
| Individual N | 4.7 | . . 2 . . . . . 2 . . 1 . . . . . . . 0 . . . |

$$\mathbf{y}_i = \sum_{j=1}^{J} x_{ij}\beta_j$$

SNPs with large $|\beta_j|$ are relevant

# Sparse Regression

- Lasso (Tibshirani, 1996) : Learn a sparse regression model

Regression coefficients

Covariates

---

# Sparse Regression

- Lasso (Tibshirani, 1996) : Learn a sparse regression model

Regression coefficients

Covariates

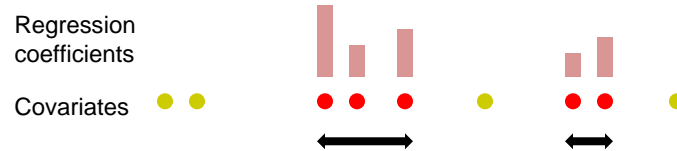- Fused lasso (Tibshirani et al., 2005) : Fuse adjacent coefficient values, assuming covariates are ordered

Regression coefficients

Covariates

# Sparse Regression

- Block-regularized regression

Regression coefficients

Covariates

- The block boundaries are determined probabilistically.

- Motivated by **association mapping** problem in computational biology
  - The block structure in genome arising from a non-random recombination process

# Recombination

| Parent chromosomes | Offspring chromosomes | After recombination |
|---|---|---|

Mother

Father

**Recombination rate** $\rho$ : frequency of recombination per unit distance on chromosome (often, per kb)

# After Many Generations with Recombination ...

Ancestor chromosomes — Descendent chromosomes

# After Many Generations with Recombination ...

Ancestor chromosomes — Descendent chromosomes

X X

X X

X X

X X

Causal SNP

8

# Bayesian Variable Selection
**(George and McCulloch, 1993, Ishwaran and Rao, 2005)**

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \epsilon, \quad \epsilon \sim N(0, \sigma^2)$$

# Bayesian Variable Selection
**(George and McCulloch, 1993, Ishwaran and Rao, 2005)**

If $C_j = 0$ (irrelevant), $\quad \beta_j = 0$

If $C_j = 1$ (relevant), use Laplacian prior

$$\beta_j | c_j \sim \frac{1}{2(2\lambda\sigma^2)} \exp\left(-\frac{|\beta_j|}{2\lambda\sigma 2}\right)$$

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \epsilon, \quad \epsilon \sim N(0, \sigma^2)$$

9

# Bayesian Variable Selection

**(George and McCulloch, 1993, Ishwaran and Rao, 2005)**

Bernoulli prior on $C_i$ 's
$$c_j \sim \text{Bernoulli}(p)$$

If $C_j = 0$ (irrelevant), $\beta_j = 0$

If $C_j = 1$ (relevant), use Laplacian prior

$$\beta_j | c_j \sim \frac{1}{2(2\lambda\sigma^2)} \exp\left(-\frac{|\beta_j|}{2\lambda\sigma 2}\right)$$

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \epsilon, \quad \epsilon \sim N(0, \sigma^2)$$

---

# Markov Chain Prior

$$P(\mathbf{c}) = P(c_1) \prod_{j=2}^{J} P(c_j | c_{j-1})$$

# Markov Chain Prior

$$P(\mathbf{c}) = P(c_1) \prod_{j=2}^{J} P(c_j|c_{j-1})$$

- $c_j = c_{j-1}$ if
    1) the distance between the two SNPs is small, or
    2) the recombination rate between the two SNPs is small

# Markov Chain Prior

$$P(\mathbf{c}) = P(c_1) \prod_{j=2}^{J} P(c_j|c_{j-1})$$

Poisson process

$$P(c_j|c_{j-1}) = \exp(-d_j\rho_j)\ \delta(c_j, c_{j-1})$$
$$+(1 - \exp(-d_j\rho_j))\ \Pi_{c_{j-1},c_j}$$

- $\rho_j$ : Recombination rate at $j$th SNP
- $d_j$ : Distance between $j$th and ($j$-1)th SNP
- $\Pi$ : Transition probability matrix

$$\begin{pmatrix} \pi_0 & 1 - \pi_0 \\ 1 - \pi_1 & \pi_1 \end{pmatrix}$$

# Block-regularized Regression with Markov Chain Prior

Recombination rate

Distance

$\pi$  $\rho_2$  $d_2$  $\rho_J$  $d_J$

Markov chain prior on $C_j$ 's

If $C_j$ = 0 (irrelevant),  $\beta_j = 0$

If $C_j$ = 1 (relevant), use Laplacian prior

$$\beta_j | c_j \sim \frac{1}{2(2\lambda\sigma^2)} \exp\left(-\frac{|\beta_j|}{2\lambda\sigma 2}\right)$$

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \epsilon, \quad \epsilon \sim N(0, \sigma^2)$$

$c_1$  $c_2$  $\cdots$  $c_J$

$\beta_1$  $\beta_2$  $\beta_J$

$X$  $y$

$\sigma^2$

---

# Learning with MCMC

- In each iteration
  - Sample $(C_j, \beta_j)$'s

$$p(\beta_j, c_j | \boldsymbol{\beta}_{-j}, \mathbf{c}_{-j}, \mathbf{y}, \mathbf{X}, \sigma^2)$$
$$= p(\beta_j | \boldsymbol{\beta}_{-j}, \mathbf{c}, \mathbf{y}, \mathbf{X}, \sigma^2)$$
$$\cdot P(c_j | \boldsymbol{\beta}_{-j}, \mathbf{c}_{-j}, \mathbf{y}, \mathbf{X}, \sigma^2)$$

  - Sample $\pi$

  - Sample $\sigma^2$

$\pi$  $\rho_2$  $d_2$  $\rho_J$  $d_J$

$c_1$  $c_2$  $\cdots$  $c_J$

$\beta_1$  $\beta_2$  $\beta_J$

$X$  $y$

$\sigma^2$

# Experiments

- Simulation study
  - Comparison with
    - Bayesian variable selection with independent Bernoulli prior
    - Lasso
    - Ridge regression
  - Simulate covariates from *ms* (Hudson, 2002)
  - Estimate recombination rates using *PHASE* (Li and Stephens, 2004)
  - 10 relevant SNPs out of 100-250 SNPs
  - 180 individuals
  - MCMC sampling for 5000 iterations after 2000 burn-in
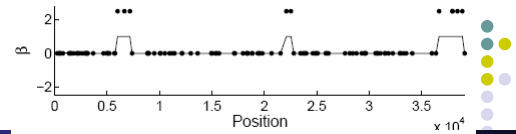
- Mouse dataset

---

## Simulations

True Model

## Slide 27

**Simulations**

True Model

Block-regularized regression

## Slide 28

**Simulations**

True Model

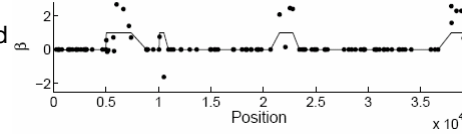Block-regularized regression

Independent Bernoulli prior

Simulations

True Model

Block-regularized regression

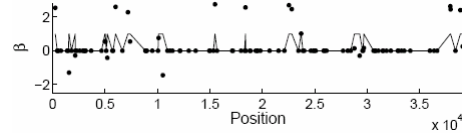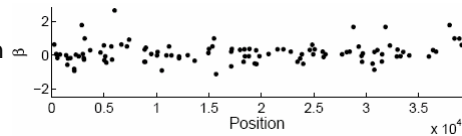Independent Bernoulli prior

Ridge regression

Simulations

True Model

Block-regularized regression

Independent Bernoulli prior
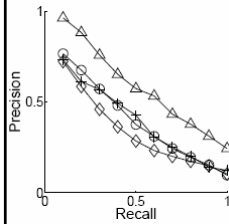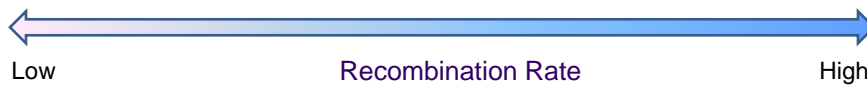
Ridge regression

Lasso

# Precision and Recall

$\rho$ = 0.05/kb

Low      Recombination Rate      High

# Precision and Recall

$\rho$ = 0.05/kb      $\rho$ = 0.1/kb      $\rho$ = 0.5/kb      $\rho$ = 1.0/kb

Low      Recombination Rate      High

# Precision and Recall



| | Legend |
|---|---|
| △ | Block |
| ⊖ | Bernoulli |
| ◇ | Ridge |
| + | Lasso |

$\rho = 0.05$/kb    $\rho = 0.1$/kb    $\rho = 0.5$/kb    $\rho = 1.0$/kb

Low ← Recombination Rate → High

---

# Mouse Data (BROAD institute)

Block-regularized regression

Independent Bernoulli prior
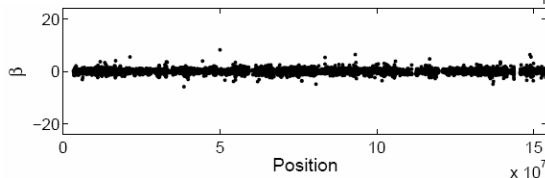
Lasso

# Multiple-trait Association

Simple (univariate) traits

**ACGTTTTACTGTACAATT**

Complex (multivariate) traits

**ACGTTTTACTGTACAATT**

$$\hat{\mathbf{B}}^{\text{lasso}} \;=\; \operatorname{argmin} \sum_{k} (\mathbf{y}_k - \mathbf{X}\boldsymbol{\beta}_k)^T \cdot (\mathbf{y}_k - \mathbf{X}\boldsymbol{\beta}_k) + \lambda \sum_{k,j} |\beta_{kj}| \quad \Rightarrow$$

Need additional constraints on **B** apart from sparsity

Eric Xing © Eric Xing @ CMU, 2006-2009 35

---

# Multiple-trait Association

**Step 1**: Thresholded correlation graph of phenotypes

body weight

blood pressure

**Overall effect of the fusion penalty**

ACGTTTTACTGTACAATT

**Step 2**: Graph-constrained fused lasso

ACGTTTTACTGTACAATT

Fusion

$$\hat{\mathbf{B}}^{\text{GC}} = \operatorname{argmin} \sum_{k} (\mathbf{y}_k - \mathbf{X}\boldsymbol{\beta}_k)^T \cdot (\mathbf{y}_k - \mathbf{X}\boldsymbol{\beta}_k)$$
$$+ \lambda \sum_{k} \sum_{j} |\beta_{jk}| + \gamma \sum_{(m,l) \in E} \sum_{j} |\beta_{jm} - \operatorname{sign}(r_{ml})\beta_{jl}|$$

Eric Xing © Eric Xing @ CMU, 2006-2009 36

18

# Multiple-trait Association
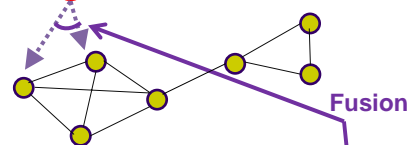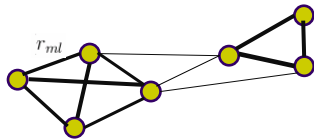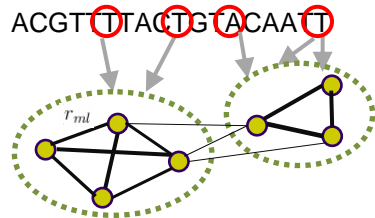
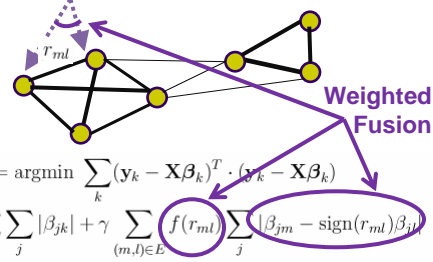**Step 1**: Thresholded correlation graph of phenotypes with **weights**

$r_{ml}$

**Overall effect of the weighted fusion penalty**

ACGTTTTACTGTACAATT

$r_{ml}$

**Step 2**: Graph-weighted fused lasso

ACGTTTTACTGTACAATT

$r_{ml}$

**Weighted Fusion**

$$\hat{\mathbf{B}}^{\mathrm{GW}} = \arg\min \sum_k (\mathbf{y}_k - \mathbf{X}\boldsymbol{\beta}_k)^T \cdot (\mathbf{y}_k - \mathbf{X}\boldsymbol{\beta}_k)$$
$$+\lambda \sum_k \sum_j |\beta_{jk}| + \gamma \sum_{(m,l)\in E} f(r_{ml}) \sum_j |\beta_{jm} - \mathrm{sign}(r_{ml})\beta_{jl}|$$

---

# Convex Optimization

- Quadratic programming formulation
  - Graph-constrained fused lasso

$$\hat{\mathbf{B}}^{\mathrm{GC}} = \arg\min \sum_k (\mathbf{y}_k - \mathbf{X}\boldsymbol{\beta}_k)^T \cdot (\mathbf{y}_k - \mathbf{X}\boldsymbol{\beta}_k)$$

  - Gra   s. t.   $\sum_k \sum_j |\beta_{jk}| \le s_1$ and $\sum_{(m,l)\in E} \sum_j |\beta_{jm} - \mathrm{sign}(r_{ml})\beta_{jl}| \le s_2$

$$\hat{\mathbf{B}}^{\mathrm{GW}} = \arg\min \sum_k (\mathbf{y}_k - \mathbf{X}\boldsymbol{\beta}_k)^T \cdot (\mathbf{y}_k - \mathbf{X}\boldsymbol{\beta}_k)$$

- Many   s. t.   $\sum_k \sum_j |\beta_{jk}| \le s_1$ and $\sum_{(m,l)\in E} f(r_{ml}) \sum_j |\beta_{jm} - \mathrm{sign}(r_{ml})\beta_{jl}| \le s_2$   g convex optimization problems can be used

# Simulated Data

Phenotype Correlation Structure



A         B

C        D        E

Regression Coefficients

F (Lasso)    G ($G_c$Flasso)    H    ($G_w$Flasso)  I

A: The correlation matrix of phenotypes, B: the edges of the phenotype correlation graph obtained at threshold 0.3 are shown as white pixels, C: The true regression coefficients used in simulation. Rows correspond to SNPs and columns to phenotypes. D: -log(p-value). Absolute values of the estimated regression coefficients are shown for E: ridge regression, F: lasso, G: $G_c$Flasso, H: $G_1^w$Flasso, and I: $G_2^w$Flasso.

---

# ROC on Simulated Data



The same trend is observed under:

- Different sample sizes,
- Association strength
- Edge weight cut-offs

20

## Asthma Multiple-trait Association



| Phenotype Correlation Structure | Single-marker Single-trait test | Lasso | Graph-constrained Fused lasso | Graph-weighted Fused lasso |

## Summary

- Likelihood-based Structured Input or Structured Output
  - Block-regularized regression makes use of the prior knowledge on the block structure such as distance and recombination rate between adjacent SNPs.
  - Graph-guided fused lasso framework incorporates correlation information among phenotypes to detect pleiotropic effect of genotypic variations.

- Future Work
  - Combine structural information in both genome and phenome in a single statistical method

# Margin-Based Discriminative Learning Paradigms

**SVM**

$$y = \text{sign}(\mathbf{w}^\top \mathbf{x} + b)$$

$$\min_{\mathbf{w},\xi} \quad \frac{1}{2}\|\mathbf{w}\|^2 + C\sum_{i=1}^{m} \xi_i$$
$$y^i(\mathbf{w}^\top \mathbf{x}^i + b) \geq 1 - \xi_i, \quad \forall i$$

**M³N**

$$y = \arg\max_{\mathbf{y}\in\mathcal{Y}(\mathbf{x})} F(\mathbf{x},\mathbf{y};\mathbf{w})$$

$$\min_{\mathbf{w},\xi} \quad \frac{1}{2}\|\mathbf{w}\|^2 + C\sum_{i=1}^{m} \xi_i$$
$$\mathbf{w}^\top[\mathbf{f}(\mathbf{x}^i) - \mathbf{f}(\mathbf{x}^i,\mathbf{y})] \geq \ell(\mathbf{y}^i,\mathbf{y}) - \xi_i, \quad \forall i, \forall \mathbf{y} \neq \mathbf{y}^i$$

**MED**

$$y = \text{sign}(\langle f(\mathbf{x},\mathbf{w})\rangle_{Q(\mathbf{w})})$$
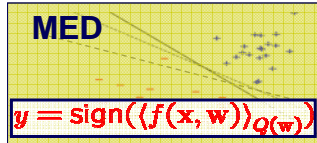
$$\min_{Q} \quad \text{KL}(Q\|Q_0)$$
$$y^i\langle f(\mathbf{x}^i)\rangle_Q \geq \xi_i, \quad \forall i$$

**MED-MN**

= SMED + Bayesian M³N

Eric Xing

© Eric Xing @ CMU, 2006-2009          43

---

# Maximum Entropy Discrimination Markov Networks

- Structured MaxEnt Discrimination (SMED):

$$P1: \quad \min_{p(\mathbf{w}),\xi} \quad KL(p(\mathbf{w})\|p_0(\mathbf{w})) + U(\xi)$$
$$\text{s.t.} \quad p(\mathbf{w}) \in \mathcal{F}_1, \ \xi_i \geq 0, \forall i.$$

*generalized* maximum entropy or *regularized* KL-divergence

- Feasible subspace of weight distribution:

$$\mathcal{F}_1 = \left\{ p(\mathbf{w}) : \int p(\mathbf{w})[\Delta F_i(\mathbf{y};\mathbf{w}) - \Delta \ell_i(\mathbf{y})]\,d\mathbf{w} \geq -\xi_i, \ \forall i, \forall \mathbf{y} \neq \mathbf{y}^i \right\},$$

*expected* margin constraints.

- Average from distribution of M³Ns

$$h_1(\mathbf{x}; p(\mathbf{w})) = \arg\max_{\mathbf{y}\in\mathcal{Y}(\mathbf{x})} \int p(\mathbf{w}) F(\mathbf{x},\mathbf{y};\mathbf{w})\,d\mathbf{w}$$

$$D(p, p_0) = KL(p\|p_0)$$

Eric Xing          © Eric Xing @ CMU, 2006-2009          44

22

# Solution to MaxEnDNet

- Theorem:
  - Posterior Distribution:

$$p(\mathbf{w}) = \frac{1}{Z(\alpha)} p_0(\mathbf{w}) \exp\left\{ \sum_{i,\mathbf{y}} \alpha_i(\mathbf{y})[\Delta F_i(\mathbf{y};\mathbf{w}) - \Delta\ell_i(\mathbf{y})] \right\}$$

  - Dual Optimization Problem:

$$\text{D1}: \quad \max_{\alpha} \; -\log Z(\alpha) - U^\star(\alpha)$$
$$\text{s.t. } \; \alpha_i(\mathbf{y}) \geq 0, \; \forall i, \; \forall \mathbf{y},$$

$U^\star(\cdot)$ is the conjugate of the $U(\cdot)$, i.e., $U^\star(\alpha) = \sup_{\xi}\left( \sum_{i,y} \alpha_i(y)\xi_i - U(\xi) \right)$

---

# Gaussian MaxEnDNet (reduction to M³N)

- Theorem
  - Assume

$$F(\mathbf{x},\mathbf{y};\mathbf{w}) = \mathbf{w}^\top \mathbf{f}(\mathbf{x},\mathbf{y}), U(\xi) = C\sum_i \xi_i, \text{and } p_0(\mathbf{w}) = \mathcal{N}(\mathbf{w}|0, I)$$

  - Posterior distribution:

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mu_{\mathbf{w}}, I), \text{ where } \mu_{\mathbf{w}} = \sum_{i,\mathbf{y}} \alpha_i(\mathbf{y})\Delta\mathbf{f}_i(\mathbf{y})$$

  - Dual optimization:

$$\max_{\alpha} \; \sum_{i,\mathbf{y}} \alpha_i(\mathbf{y})\Delta\ell_i(\mathbf{y}) - \frac{1}{2}\|\sum_{i,\mathbf{y}} \alpha_i(\mathbf{y})\Delta\mathbf{f}_i(\mathbf{y})\|^2$$
$$\text{s.t. } \sum_{\mathbf{y}} \alpha_i(\mathbf{y}) = C; \; \alpha_i(\mathbf{y}) \geq 0, \; \forall i, \; \forall \mathbf{y},$$

  - Predictive rule:

$$h_1(\mathbf{x}) = \arg\max_{\mathbf{y}\in\mathcal{Y}(\mathbf{x})} \int p(\mathbf{w})F(\mathbf{x},\mathbf{y};\mathbf{w})\,d\mathbf{w} = \arg\max_{\mathbf{y}\in\mathcal{Y}(\mathbf{x})} \mu_{\mathbf{w}}^\top \mathbf{f}(\mathbf{x},\mathbf{y})$$

- Thus, MaxEnDNet subsumes M³Ns and admits all the merits of max-margin learning
- Furthermore, MaxEnDNet has at least three advantages …

# Three Advantages

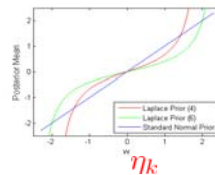- An averaging Model: PAC-Bayesian prediction error guarantee

$$\Pr_Q(M(h, \mathbf{x}, \mathbf{y}) \le 0) \le \Pr_{\mathcal{D}}(M(h, \mathbf{x}, \mathbf{y}) \le \gamma) + O\left(\sqrt{\frac{\gamma^{-2}KL(p\|p_0)\ln(N|\mathcal{Y}|) + \ln N + \ln \delta^{-1}}{N}}\right).$$

- Entropy regularization: Introducing useful biases
    - Standard Normal prior => reduction to standard M³N (we've seen it)
    - Laplace prior => Posterior shrinkage effects (sparse M³N)

$$\forall k, \ \langle w_k \rangle_p = \frac{2\eta_k}{\lambda - \eta_k^2}$$

$\eta_k$

- Integration of Generative and Discriminative principles
    - Incorporate latent variables and structures (PoMEN)
    - Semisupervised learning (with partially labeled data)

---

# Key Challenges

- Extremely high dimensionality and low data volume
    - d ~ 1M
    - N ~ 1K
    - Sample complexity with bounded error?

- Sparsity bias of the model
    - Often <100 features out of the !M are relevant
    - Regularization schemes to enforce sparsity

- Structures and hidden variables
    - Inputs and outputs often bear intricate structures (e.g., chain or graphical dependencies)
    - How to capture other latent structures between unobserved variables

- Generalizability and scalability
    - Move efficient convex opt solver and Bayesian inference algorithms

- Provable theoretical guarantees
    - Consistency and sparsistency
    - Stability, convergence rate, etc.