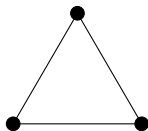


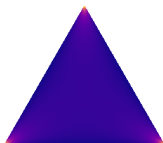
Априорное распределение на структуре модели

Каждая точка на симплексе задает модель.

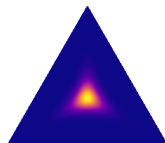
Распределение Гумбель-софтмакс: $\Gamma \sim \text{GS}(\mathbf{s}, \lambda_{\text{temp}})$



$$\lambda_{\text{temp}} \rightarrow 0$$

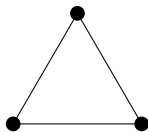


$$\lambda_{\text{temp}} = 0.995$$

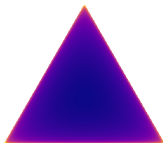


$$\lambda_{\text{temp}} = 5.0$$

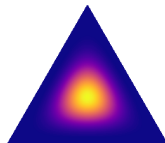
Распределение Дирихле: $\Gamma \sim \text{Dir}(\mathbf{s}, \lambda_{\text{temp}})$



$$\lambda_{\text{temp}} \rightarrow 0$$



$$\lambda_{\text{temp}} = 0.995$$

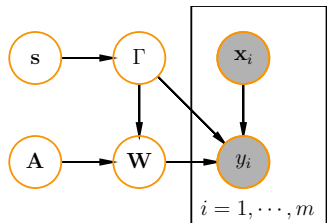


$$\lambda_{\text{temp}} = 5.0$$

Байесовский выбор модели

Базовая модель:

- параметры модели $\mathbf{w} \sim \mathcal{N}(0, \alpha^{-1})$,
- гиперпараметры модели $\mathbf{h} = [\alpha]$.



Предлагаемая модель:

- параметры модели $\mathbf{w}_r^{j,k} \sim \mathcal{N}(0, \gamma_r^{j,k} (\mathbf{A}_r^{j,k})^{-1})$, $\mathbf{A}_r^{j,k}$ — диагональная матрица параметров, соответствующих базовых функций $\mathbf{g}_r^{j,k}$, $(\mathbf{A}_r^{j,k})^{-1} \sim \text{inv-gamma}(\lambda_1, \lambda_2)$,
- структурные параметры модели $\Gamma = \{\gamma^{j,k}, (j, k) \in E\}$, $\gamma^{j,k} \sim \text{GS}(s^{j,k}, \lambda_{\text{temp}})$,
- гиперпараметры модели $\mathbf{h} = [\text{diag}(\mathbf{A}), s]$,
- метапараметры $\lambda_1, \lambda_2, \lambda_{\text{temp}}$.

Верхняя оценка правдоподобия модели

Оптимальные факторы $q^*(\theta)$, $q^*(\mathbf{m}, \mathbf{V}, \alpha)$ имеют вид

$$\ln q^* = E[\ln p(\mathbf{Z}, \theta, \mathbf{m}, \mathbf{V}, \alpha)] + \text{const}$$

и не вычисляются аналитически, так как правдоподобие L в модели $p(\mathbf{Z}, \theta, \mathbf{m}, \mathbf{V}, \alpha)$ содержит сумму экспонент $g(\mathbf{s}_n)$ в знаменателе,

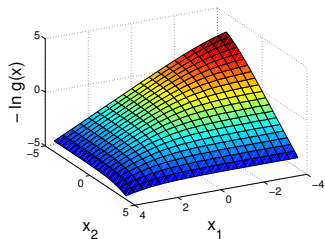
$$L(\mathbf{Z}|\theta, \alpha) = \prod_{n=1}^N \prod_{k=1}^{K_h} \left[\frac{\exp(s_{n,k})}{g(\mathbf{s}_n)} \right]^{z_{nk}}, \quad g(\mathbf{s}_n) = \sum_{k=1}^{K_h} \exp(s_{n,k}).$$

$\ln \frac{1}{g(\mathbf{x})}$ – вогнутая функция. Касательная плоскость к ней через точку ξ :

$$\zeta(\mathbf{x}, \xi) = -\ln(g(\xi)) - \nabla \ln(g(\xi))^T (\mathbf{x} - \xi).$$

Верхняя оценка правдоподобия L

$$\frac{1}{g(\xi)} \exp \left(s_{n,k} + \sum_{k'=1}^{K_h} \frac{\exp(\xi_{k'})}{g(\xi)} (\xi_{k'} - s_{n,k'}) \right).$$



Анализ качества предлагаемого метода выбора признаков

Решена задача прогнозирования многомерных временных рядов $\mathbf{y}(t) \in \mathbb{R}^3$ координат конечности по интервалам $\mathbf{s}(t - \Delta t)$ многомерных временных рядов $\mathbf{s}(t) \in \mathbb{R}^{N_{\text{ch}}}$ многоканальных электрокортикограмм.

Признаки:

$$\underline{\mathbf{X}}_m \in \mathbb{R}^{F \times N_{\text{ch}}}, \quad \underline{\mathbf{X}}_{mjn} = \begin{cases} s_n(t_m + \tau), & j = 1, \\ W_{mjn}, & j = 2, \dots, F + 1, \end{cases}$$

Прогноз в точке t_m :

$$\hat{\mathbf{y}}_m = \text{vec}(\underline{\mathbf{X}}_m)^T \hat{\mathbf{w}}.$$

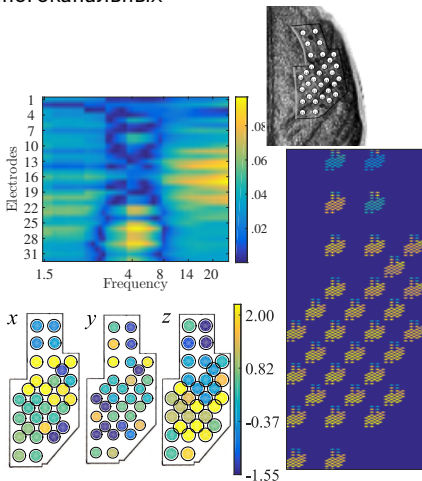
Качество прогнозирования:

коэффициент корреляции между $\hat{\mathbf{Y}}$ и \mathbf{Y} ,

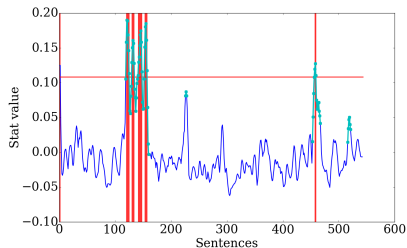
$$\text{corr}(\hat{\mathbf{Y}}, \mathbf{Y}) = \frac{\text{cov}(\hat{\mathbf{y}}, \mathbf{y})}{\sqrt{\text{cov}(\hat{\mathbf{y}}, \hat{\mathbf{y}})\text{cov}(\mathbf{y}, \mathbf{y})}}.$$

масштабированная ошибка MSE,

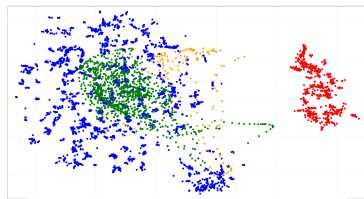
$$\text{sMSE}(\hat{\mathbf{Y}}, \mathbf{Y}) = \frac{\sum_{m=1}^M \|\hat{\mathbf{y}}_m - \mathbf{y}_m\|_2}{\sum_{m=1}^M \|\bar{\mathbf{y}} - \mathbf{y}_m\|_2}.$$



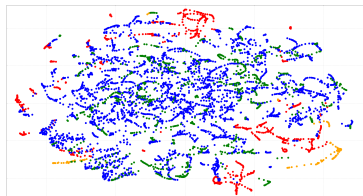
Снижение размерности с сохранением локальной структуры близости



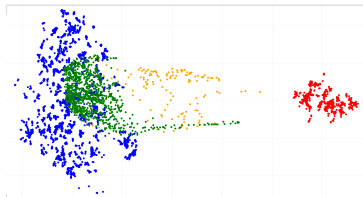
Построение признакового пространства для ряда s



Модифицированный t-SNE, $\mu = 10$



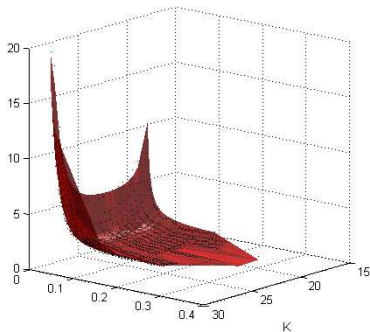
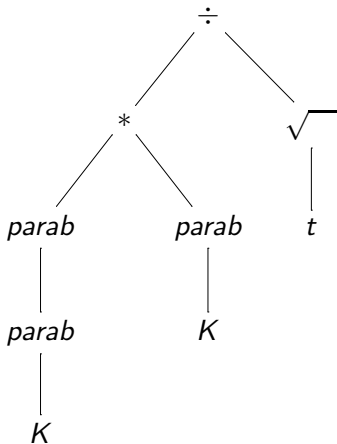
Исходный t-SNE



Модифицированный t-SNE, $\mu = 100$

Результаты для опциона CLZ

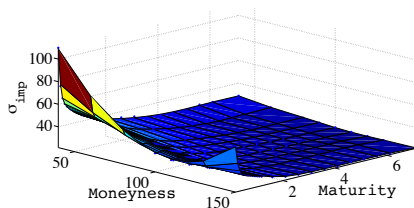
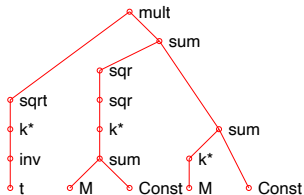
$$\sigma = \frac{(w_1 K + w_2)(w_1 K^2 + w_2 K + w_3)^2}{\sqrt{t}}$$



Resulting models

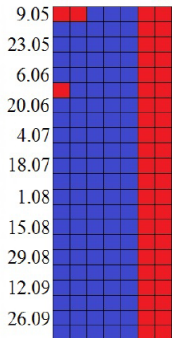
Resulting model

$$\sigma_{imp} = \frac{(k_2 M + k_3)^4 + k_4 M + c}{\sqrt{k_1 t}}$$

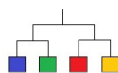
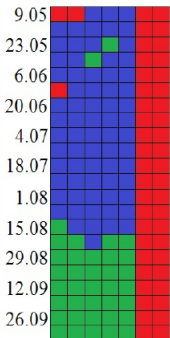




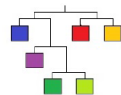
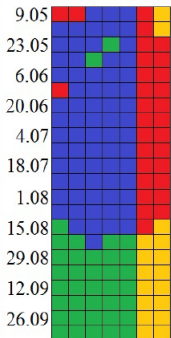
Пн Bc



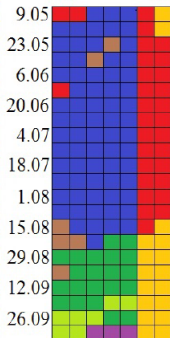
Пн Bc

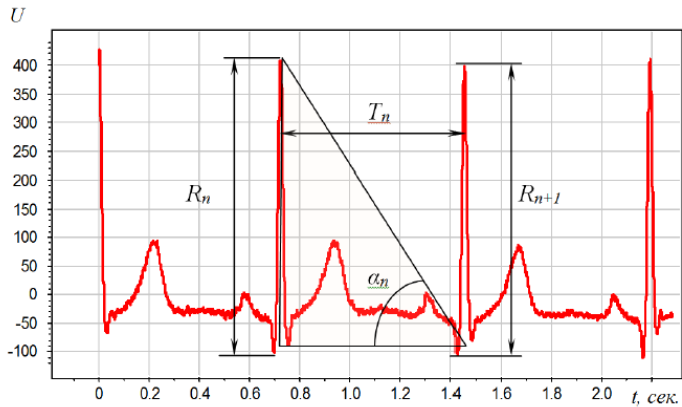


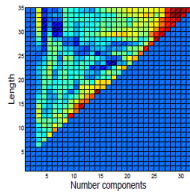
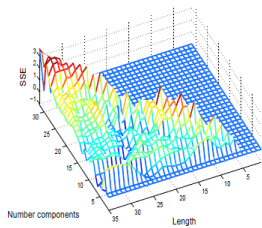
Пн Bc

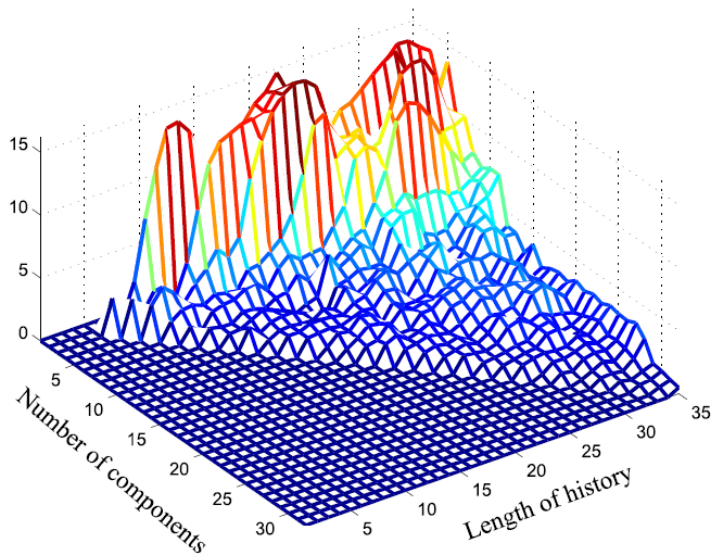


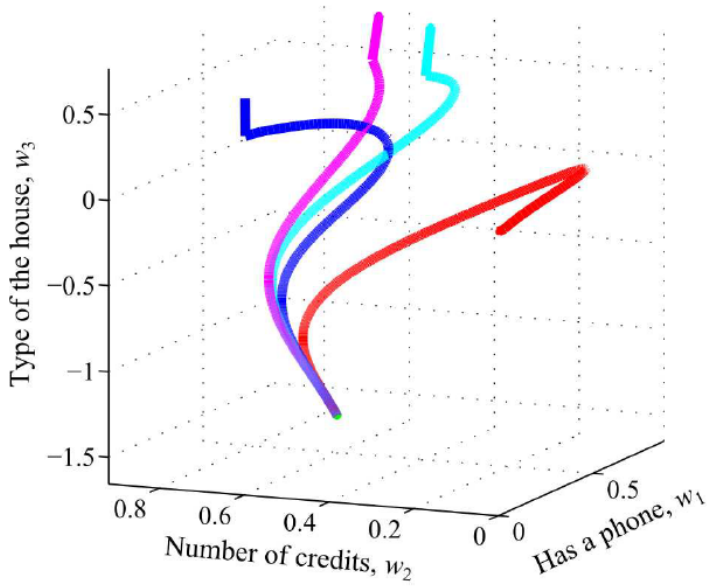
Пн Bc

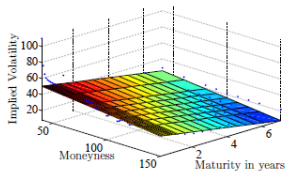
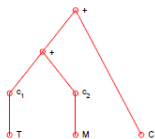




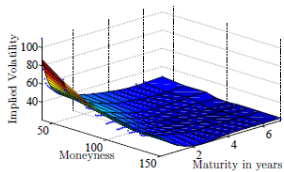
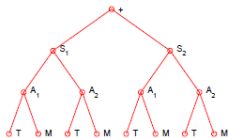




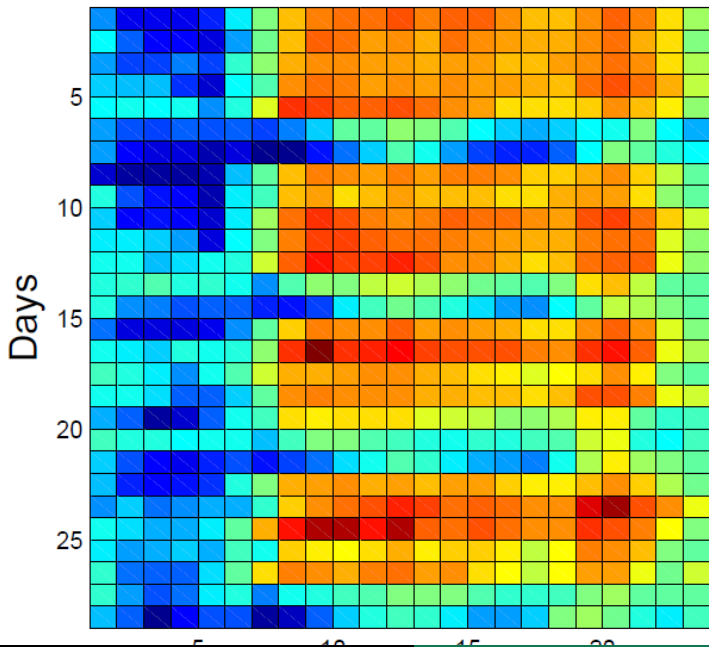


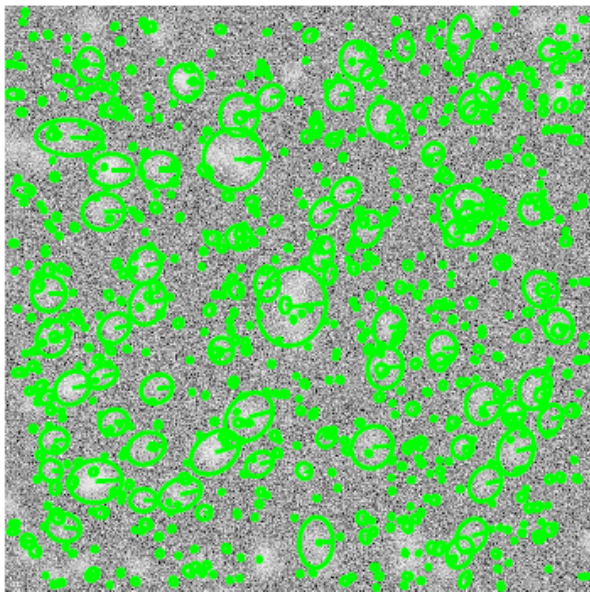


a)

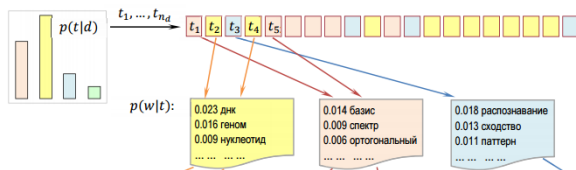


b)





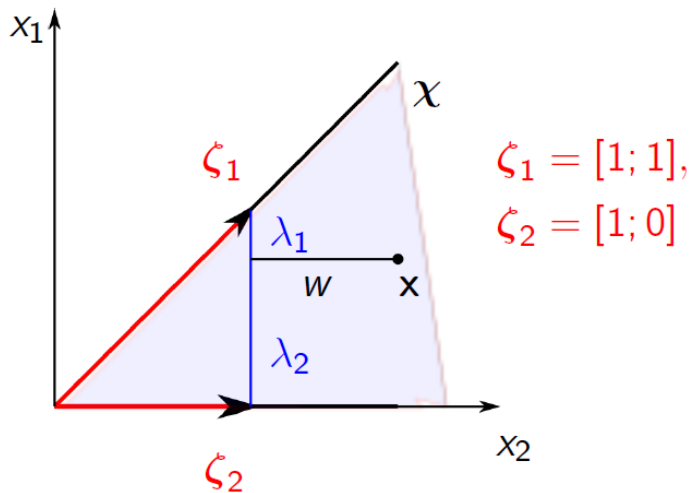
Тематическое моделирование и матричное разложение

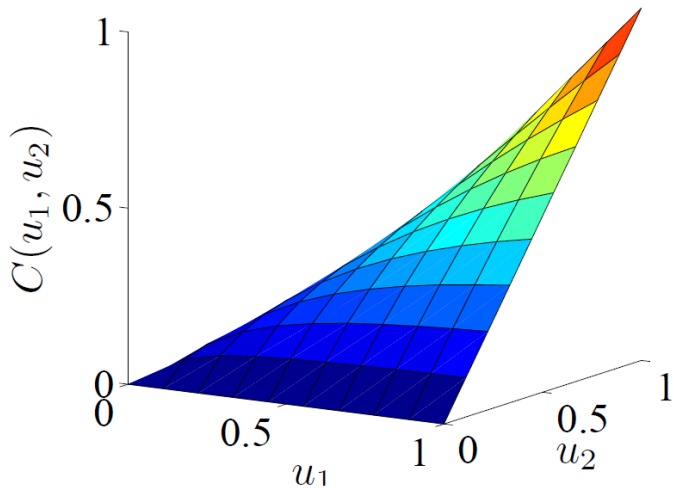


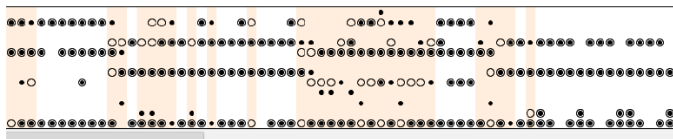
w_1, \dots, w_{n_d} :

Разработан спектрально-аналитический подход к выявлению размытых протяженных повторов в геномных последовательностях. Метод основан на разномасштабном оценивании сходства нуклеотидных последовательностей в пространстве коэффициентов разложения фрагментов кривых GC- и GA-содержания по классическим ортогональным базисам. Найдены условия оптимальной аппроксимации, обеспечивающие автоматическое распознавание повторов различных видов (прямых и инвертированных, а также тандемных) на спектральной матрице сходства. Метод одинаково хорошо работает на разных масштабах данных. Он позволяет выявлять следы сегментных дупликаций и мегасателлитные участки в геноме, районы синтении при сравнении пары геномов. Его можно использовать для детального изучения фрагментов хромосом (поиска размытых участков с умеренной длиной повторяющегося паттерна).

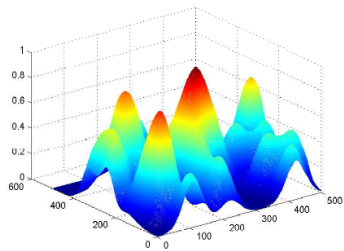
Vorontsov K. V., Potapenko A. A., Plavin A. V. Additive Regularization of Topic Models for Topic Selection and Sparse Factorization // International Symposium On Learning And Data Sciences 2015.



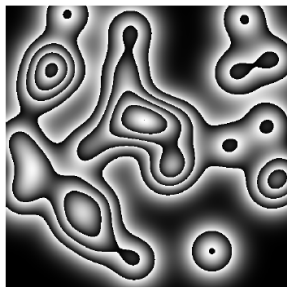




Пустые круги — истинные полутона, точки — предсказанные, ошибки подсвечены, горизонтальная ось — время.



(a) Высота



(б) Фазовая составляющая без учета шума

