

Survival driven deconvolution (DeSurv) reveals prognostic and interpretable cancer subtypes

Amber M. Young^{a,1,2}, Alisa Yurovsky^b, Didong Li^a, and Naim U. Rashid^{a,c}

^aUniversity of North Carolina at Chapel Hill, Biostatistics, Street, City, State, Zip; ^bStony Brook University, Street, City, State, Zip

This manuscript was compiled on January 22, 2026

Molecular subtyping in cancer is an ongoing problem that relies on the identification of robust and replicable gene signatures. While transcriptomic profiling has revealed recurrent gene expression patterns in various types of cancer, the prognostic value of these signatures is typically evaluated in retrospect. This is due to the reliance on unsupervised learning methods for identifying cell-type-specific signals and clustering patients into molecular subtypes. Here we present a Survival-driven Deconvolution tool (deSurv) that integrates bulk RNA-sequencing data with patient survival information to identify cell-type-enriched gene signatures associated with prognosis. Applying deSurv to various cohorts in pancreatic cancer, we uncover prognostic and biologically interpretable subtypes that reflect the complex interactions between stroma, tumor, and immune cells in the tumor microenvironment. Our approach highlights the value of using patient outcomes during gene signature discovery.

one | two | optional | optional | optional

Molecular subtyping has transformed precision oncology by stratifying patients into biologically and clinically meaningful groups that inform prognosis and guide therapy (1–5). Subtyping relies on the identification of robust biological signals that define subtypes such as transcriptomic signatures. However, the tumor microenvironment (TME) contains mixtures of diverse cell types such as malignant, stromal, immune, and endothelial cells, and disentangling tumor specific signals from this mixture can be challenging. As such, subtyping pipelines typically rely on the deconvolution of bulk transcriptomic data or single-cell analysis to discover distinct cell types and their corresponding signatures. Downstream, the signatures are evaluated for clinical relevance such as overall survival or response to treatment.

Separating discovery from validation can risk overfitting and limits biological and clinical generalizability. Identified cell types may capture dataset-specific noise rather than reproducible biological signals, undermining their utility in downstream analyses or therapeutic targeting (6, 7). Moreover, even when discovered cell types are biologically valid and reproducible, they may not correspond to the cellular programs most relevant for predicting or influencing clinical outcomes (8, 9). Therefore, there is a clear need for integrative methods that jointly uncover biologically meaningful programs while directly incorporating clinical endpoints to ensure prognostic relevance.

However, integrating patient outcomes into the discovery phase is not straightforward with current technology and methodology. Single-cell transcriptomics can resolve programs at the cellular level, but cohort sizes are often too small to support survival analyses. In contrast, large bulk transcriptomic cohorts with clinical annotations are well-suited for outcome modeling (10, 11), yet deconvolution is needed to disentangle overlapping cellular signals. Reference-based deconvolution

methods focus on estimating cell-type proportions from predefined signatures, which limits the utility of these methods for discovery of novel programs (12).

Nonnegative matrix factorization (NMF) is widely used in cancer genomics because its nonnegativity constraints produce biologically interpretable, additive molecular programs (13–16). Although recent extensions have incorporated supervision into the factorization, most target regression or classification rather than time-to-event outcomes. Two studies have proposed survival-aware NMF formulations (17, 18), but both integrate the survival objective through the sample-specific loadings rather than the gene-level programs. This design emphasizes prediction accuracy but limits the model's ability to restructure or refine the underlying molecular programs, reducing its value for biological interpretation and subtype discovery, which are core objectives in cancer transcriptomics. In addition, neither study provides a principled approach for hyperparameter selection or model assessment, and convergence properties are only briefly addressed in one manuscript. Both works remain unpublished and unreviewed, leaving their methodological robustness and reproducibility uncertain. These gaps highlight the need for a rigorously formulated, survival-aware deconvolution method that jointly estimates interpretable molecular programs and their prognostic relevance.

Here we present DeSurv, a Survival-supervised Deconvolution framework that integrates non-negative matrix factor-

Significance Statement

Tumor transcriptomes reflect mixtures of malignant and microenvironmental cell populations, making it challenging to identify the molecular programs that truly drive clinical outcomes. Existing deconvolution and matrix factorization methods discover latent transcriptional programs but do not ensure that these programs are prognostic, while supervised extensions optimized for prediction offer limited biological interpretability. We present DeSurv, a survival-supervised deconvolution framework that integrates nonnegative matrix factorization with Cox modeling to jointly learn biologically coherent gene programs and their associations with patient survival. By embedding outcome information directly into the discovery process and performing automatic model selection, DeSurv reveals clinically relevant transcriptional programs that are reproducible across cohorts. This approach advances the statistical foundations of tumor deconvolution and provides a general tool for identifying actionable molecular drivers of disease progression.

Please provide details of author contributions here.

Please declare any conflict of interest here.

² To whom correspondence should be addressed. E-mail: ayoung31@live.unc.edu

ization (NMF) with Cox proportional hazards modelling. In contrast to fully unsupervised approaches that evaluate survival associations only after the factorization, and to existing supervised NMF models that link outcomes to the subject-level factor loadings, DeSurv integrates survival information directly into the gene signature matrix. This design ensures that the discovered transcriptional programs are not only biologically interpretable but also intrinsically aligned with patient outcomes. To enhance robustness and reproducibility, DeSurv performs automatic parameter selection via Bayesian optimization, addressing the quintessential challenge of rank determination in matrix factorization.

By coupling latent program discovery with direct survival supervision, DeSurv resolves longstanding challenges in disentangling tumor–microenvironment interactions and aligns molecular heterogeneity with clinical outcomes. This unified approach represents a methodological advance in translational cancer genomics and provides a general framework for deriving actionable insights from high-dimensional transcriptomic data.

Results

Model Overview. We have developed an integrated framework, DeSurv, that couples Nonnegative Matrix Factorization (NMF) with Cox proportional hazards regression to identify latent gene-expression programs associated with patient survival (Figure 1A). The model takes as input a bulk expression matrix of p genes by n patients (X) together with corresponding survival times (y) and censoring indicators (δ) (Figure 1A).

DeSurv optimizes a joint objective combining the NMF reconstruction loss and the Cox model’s log-partial likelihood, weighted by a supervision parameter (α) that determines the relative contribution of each term (fig. 1B):

$$(1 - \alpha) \mathcal{L}_{NMF}(X \approx WH) - \alpha \mathcal{L}_{Cox}(X^T W \beta, y, \delta) \quad [1]$$

When $\alpha = 0$, the method reduces to standard unsupervised NMF; when $\alpha > 0$, survival information directly guides the learned factors toward prognostic structure.

Within this framework, the product ($X^T W$) represents patient-level factor scores - the inferred burden of each latent program across subjects. These factor scores serve as covariates in the Cox model, and their regression coefficients (β) indicate whether higher activity of a given program corresponds to improved or reduced survival.

Model training yields gene weights (\hat{W}), factor loadings (\hat{H}), and Cox coefficients ($\hat{\beta}$) (Figure 1C), where the inner dimension (k) specifies the number of latent factors. Genes with high gene weights in one factor and low gene weights in all others define the factor-specific signature genes (Figure 1D). By integrating survival supervision into the factorization, DeSurv not only reconstructs the underlying expression structure, preserving biological interpretability, but also guides latent factors to be prognostically informative. Subsequent analyses can therefore focus on the survival-associated gene programs (Figure 1E).

Bayesian optimisation selects the DeSurv hyperparameters (k, α).

A. Outcome-guided model selection resolves ambiguity in NMF rank choice. We examined the problem of selecting the number of latent components (k) in nonnegative matrix factorization (NMF) using gene expression data from pancreatic ductal adenocarcinoma (PDAC) cohorts. These heterogeneous tumor transcriptomes provide a representative setting in which to evaluate how commonly used unsupervised rank-selection heuristics behave in practice.

Across a range of candidate ranks, standard NMF diagnostics yielded inconsistent guidance (Fig. 2A-C). Reconstruction residuals decreased smoothly with increasing k and did not exhibit a clear elbow, a pattern consistent with both relatively small solutions ($k \approx 3-4$) and substantially larger ranks ($k \approx 6-8$). The cophenetic correlation coefficient began to decline at low ranks ($k \approx 3-4$) but continued to fluctuate at higher values without a distinct transition point. In contrast, mean silhouette width, evaluated across multiple distance metrics, was highest at very small ranks ($k \approx 2-3$) and decreased monotonically thereafter, favoring low-dimensional solutions that conflicted with recommendations based on reconstruction error or cophenetic correlation. Together, these unsupervised criteria pointed to different and incompatible values of k , highlighting the ambiguity of rank selection in standard NMF when applied to PDAC data.

To resolve this ambiguity, we applied DeSurv, which incorporates survival outcomes directly into the factorization process and evaluates models using survival-based predictive performance. Using the same PDAC gene expression data, we assessed model performance across the joint space of the number of components (k) and supervision strength (α) using cross-validated concordance index (C-index). The resulting C-index surface summarizes expected predictive performance across candidate models and enables direct comparison of solutions that differ in both model complexity and degree of supervision (Fig. 2D).

Model selection was based on standard cross-validation principles. Rather than selecting the single parameter combination with the highest predicted C-index, we selected the smallest value of k whose predicted performance lay within one standard error of the maximum. This criterion yielded a stable and parsimonious choice of model rank in the PDAC data, in contrast to the conflicting recommendations produced by unsupervised NMF heuristics.

To further evaluate rank recovery under controlled conditions, we conducted simulation studies in which the true underlying rank was known ($k = 3$). Across repeated simulation replicates, DeSurv consistently selected the correct rank, producing a concentrated distribution of selected k values centered at the true value. In contrast, standard NMF followed by post hoc Cox modeling ($\alpha = 0$) exhibited substantially greater variability and a systematic tendency toward under-selection. Together, these results indicate that incorporating outcome information during model fitting improves the reliability of rank selection in settings where unsupervised criteria yield conflicting conclusions.

“(A-D) Analyses based on real pancreatic ductal adenocarcinoma (PDAC) gene expression data illustrate the ambiguity of rank selection in standard nonnegative matrix factorization (NMF) and the use of outcome supervision in DeSurv. (A-C) Commonly used unsupervised heuristics for selecting the number

of components (k) yield inconsistent conclusions. (A) Reconstruction residuals decrease smoothly with increasing k and do not exhibit a clear elbow; diminishing returns could be inferred at intermediate ($k \approx 3$ -4) or larger ($k \approx 6$ -8) ranks. (B) The cophenetic correlation coefficient, often used to select the largest k prior to a marked loss of clustering stability, begins to decline at low ranks ($k \approx 3$ -4) but continues to fluctuate thereafter, providing no unambiguous selection criterion. (C) Mean silhouette width across multiple distance metrics is highest at small ranks ($k \approx 2$ -3) and decreases monotonically with increasing k , favoring lower-dimensional solutions that conflict with the other criteria. (D) Heatmap of the Gaussian process predicted mean cross-validated concordance index (C-index) from Bayesian optimization over the joint space of the number of components (k) and supervision strength (α), computed on the same PDAC data. The predicted performance surface summarizes survival prediction accuracy across parameter settings and illustrates how DeSurv uses outcome information to inform model selection. (E) Results from simulation studies with a known underlying rank ($k = 3$) showing the distribution of selected k values across repeated replicates. DeSurv more consistently recovers the true rank, yielding a concentrated distribution centered at $k = 3$, whereas standard NMF with post hoc Cox modeling ($\alpha = 0$) exhibits greater variability and a tendency toward under-selection. \label{fig:bo}”, fig.env=’figure’, fig.pos=’t’, out.height= “5.5in”, out.width=“6in”}

```
tar_load(fig_bo_heat_tcgacptac)
tar_load(fig_residuals_tcgacptac)
tar_load(fig_cophenetic_tcgacptac)
tar_load(fig_silhouette_tcgacptac)
```

1. upper = plot_grid(fig_bo_cvk_tcgacptac,fig_bo_cvalpha_tcgacptac,labels=c(“A”,“B”))

```
upper = plot_grid(fig_residuals_tcgacptac,
fig_cophenetic_tcgacptac, fig_silhouette_tcgacptac, ncol =
3, labels = c(“A”,“B”,“C”))
k_hist = plot_grid(alt_plots$k_hist,NULL,nrow=2,rel_heights
= c(4,1))
lower = plot_grid(fig_bo_heat_tcgacptac,
k_hist,ncol=2,labels=c(“D”,“E”),rel_widths = c(3,2))
plot_grid(upper,lower,nrow=2,rel_heights = c(2,3))
```

DeSurv improves selection of prognostic gene signatures and supervised survival analysis. To test whether supervision improves the selection of prognostic gene signatures, we used simulations with known lethal

```
{\centering \includegraphics[width=\textwidth]{/work/users/a/y/a/figures/figure10.pdf}}
```

\caption{Performance comparison between DeSurv and unsupervised NMF ($\alpha = 0$) in simulation. (A) Distribution of c

Survival-informed factorization reorganizes transcriptional structure. To assess the biological structure captured by standard nonnegative matrix factorization (NMF) and DeSurv, we examined

In the standard NMF solution, factors largely reflected dominant

By contrast, DeSurv produced a factorization that emphasized axo

These differences were reflected quantitatively by contrasting t

To further characterize how DeSurv reorganizes the transcriptional structure, we examined the gene-level structure captured by DeSurv. \begin{figure*}[t]

```
{\centering \includegraphics[width=\textwidth,height=4.5in]{/work/users/a/y/a/figures/figure11.pdf}}
```

\caption{Survival-informed factorization reorganizes transcriptional structure. To assess the biological structure captured by standard nonnegative matrix factorization (NMF) and DeSurv, we examined

DeSurv-derived latent structure generalizes to independent data. To assess generalization of DeSurv latent factors, the gene-level structure captured by DeSurv was compared to the gene-level structure captured by standard NMF.

For each method, we focused on the factor showing the largest in

When validation samples were pooled and stratified into high- and low-risk groups, the gene-level structure captured by DeSurv was

Together, these results indicate that DeSurv identifies latent factors that are biologically meaningful. \begin{figure*}[t]

```
{\centering \includegraphics[width=\textwidth]{/work/users/a/y/a/figures/figure12.pdf}}
```

\caption{Figure X. DeSurv learns prognostic structure that generalizes to independent data. To assess generalization of DeSurv latent factors, the gene-level structure captured by DeSurv was compared to the gene-level structure captured by standard NMF.

Bladder cancer analysis {.unnumbered}

We will report a focused bladder cancer analysis here, including

Discussion {.unnumbered}

We present DeSurv, a survival-driven deconvolution framework that

Results {.unnumbered}

In the standard NMF solution, factors largely reflected dominant

By contrast, DeSurv produced a factorization that emphasized axo

These differences were reflected quantitatively by contrasting t

Materials and methods {.unnumbered}

Problem formulation and notation

Let $\mathbf{X} \in \mathbb{R}^{p \times n}$, $\mathbf{Y} \in \mathbb{R}^{p \times 1}$ denote the nonnegative

The DeSurv Model

DeSurv improves selection of prognostic gene signatures and supervised survival analysis. To test whether supervision improves the selection of prognostic gene signatures, we used simulations with known lethal

```

The joint objective is
\begin{equation}
\label{eqn:desurv}
\mathcal{L}(W,H,\beta) =
(1-\alpha)\mathcal{L}_{\mathrm{NMF}}(W,H)
- \alpha\mathcal{L}_{\mathrm{Cox}}(W,\beta),
\end{equation}
where  $\mathcal{L}_{\mathrm{NMF}}(W,H)$  is the NMF reconstruction loss,  $\mathcal{L}_{\mathrm{Cox}}(W,\beta)$  is the Cox loss, and
 $\beta = (\beta_1, \dots, \beta_p)$  is the vector of Cox coefficients.

## Hyperparameter selection and cross-validation
Hyperparameters  $(k, \alpha, \lambda_H, \lambda, \xi)$  were selected by maximizing the cross-validated DeSurv loss.

## Simulation studies
Simulation studies were conducted to assess recovery of parameters and survival prediction. Gene expression data were generated from an exponential proportional hazards model, and survival times were generated from an exponential distribution.

## Evaluation metrics
where  $\|\cdot\|_F$  is the Frobenius norm and  $\beta_k$  are the Cox coefficients. The Cox log-likelihood is defined as
\begin{equation}
\ell(W, \beta) = -\sum_{i=1}^n \log \left( \sum_{j=1}^p \exp(\beta_j X_{ij}) \right) - \sum_{i=1}^n \beta_i X_{i1}
\end{equation}
Hyperparameters satisfy  $\alpha \in [0,1]$ ,  $\lambda_H \geq 0$ ,  $\lambda \geq 0$ , and  $\xi \in [0,1]$ . The constants  $1/(2np)$ ,  $2/n_{\text{event}}$ , and  $1/(2nk)$  are for numerical convenience and do not affect the minimizer.

\subsection{Optimization Algorithm}

\subsubsection{Block coordinate descent scheme}
Algorithm~\ref{alg:desurv} summarizes the block coordinate descent scheme used to minimize  $\mathcal{L}(W,H,\beta)$ . At each outer iteration we update  $H$ , then  $W$ , then  $\beta$  while holding the other blocks fixed.

\begin{algorithm}[H]
\caption{DeSurv block coordinate descent}
\label{alg:desurv}
\begin{algorithmic}[1]
\Require  $X \in \mathbb{R}^{\geq 0}_{p \times n}$ , survival times  $\tau \in \mathbb{R}^n$ 
\Ensure Fitted DeSurv parameters  $(W,H,\beta)$ 

\State Initialize  $W^{(0)}, H^{(0)}$  with positive entries
\State Initialize  $\beta^{(0)}$  (e.g.,  $\beta^{(0)} = \mathbf{0}$ )
\State  $\text{loss} = \mathcal{L}(W^{(0)}, H^{(0)}, \beta^{(0)})$ 


```

```

\State $seps = \infty$, $t = 0$
\While{$seps \geq \text{tol}$ \textbf{and} $t < \text{maxit}$}{
  \State \textbf{(H-update)}
  \State \hspace{0.5cm} Update  $H^{(t+1)}$  from  $H^{(t)}$  using the multiplicative rule in Eq.~\ref{eqn:Hupda}
  \State \hspace{0.5cm} holding  $W^{(t)}$  and  $\beta^{(t)}$  fixed.
  \State \textbf{(W-update)}
  \State \hspace{0.5cm} Update  $W^{(t+1)}$  from  $W^{(t)}$  using the hybrid multiplicative rule in Eq.~\ref{eqn:Wupdate}
  \State \hspace{0.5cm} holding  $H^{(t+1)}$  and  $\beta^{(t+1)}$  fixed, and perform backtracking to ensure
  \State \hspace{0.5cm}  $\mathcal{L}$  does not increase.
  \State \textbf{(\beta-update)}
  \State \hspace{0.5cm} Update  $\beta^{(t+1)}$  from  $\beta^{(t)}$  using a Newton-like step for the Cox loss
  \State \hspace{0.5cm} (Eq.~\ref{eqn:beta_update}), holding  $W^{(t+1)}$  and  $H^{(t+1)}$  fixed.
  \State  $\text{lossNew} = \mathcal{L}(W^{(t+1)}, H^{(t+1)}, \beta^{(t+1)})$ 
  \State  $seps = |\text{lossNew} - \text{loss}| / |\text{loss}|$ 
  \State  $\text{loss} = \text{lossNew}$ 
  \State  $t = t + 1$ 
}
\EndWhile
\State \Return  $\hat{W} = W^{(t+1)}$ ,  $\hat{H} = H^{(t+1)}$ ,  $\hat{\beta} = \beta^{(t+1)}$ 
\end{algorithmic}
\end{algorithm}

\subsubsection{Update for  $W$ }
Conditional on  $W$ , the loss  $\mathcal{L}(W, H, \beta)$  reduces to a convex quadratic function of  $W$ . We adopt the standard NMF multiplicative
update with  $\ell_2$  penalty:
\begin{equation}
H \leftarrow \max\left( \frac{W^{\top} X}{W^{\top} W H + \lambda_H H + \epsilon_H}, \epsilon_H \right),
\end{equation}
where  $\odot$  denotes elementwise multiplication, all divisions are
elementwise, and  $\epsilon_H > 0$  is a small floor to prevent the denominator
becoming exactly zero. This update is equivalent to a majorization step and
guarantees a nonincreasing reconstruction term conditional on  $W$ 
[seung2001algorithms; pascualmontano2006nonsmooth]. The same steps applied to  $W$ 
ensures that iterates remain in the interior of the nonnegative orthant, which
simplifies the convergence analysis.

\subsubsection{Update for  $W$ }
For  $W$ , we construct a hybrid multiplicative update that combines the
contributions of the NMF and Cox gradients. Let

$$\nabla_W \mathcal{L}_{\text{NMF}}(W, H) = \frac{1}{n} (W H - X) H^{\top},$$

and let

$$\nabla_W \mathcal{L}_{\text{Cox}}(W, \beta) = \frac{2}{n_{\text{event}}} \nabla_W \ell(W, \beta),$$

where  $\nabla_W \ell(W, \beta)$  denotes the Cox gradient with respect to  $W$ . Because a
nonincrease in the full loss  $\mathcal{L}$ , we embed it in a back
line search.

```

```

\begin{algorithm}[H]
\caption{$W$ update with backtracking}
\label{alg:backtrack}
\begin{algorithmic}[1]
\Require Value of $W$ and the previous iteration $W^{(t)}$, the multiplicative update term at the current
\Ensure $W^{(t+1)}$
\State $\theta = 1$
\State $b = 1$
\State $\text{flag\_accept} = \text{FALSE}$
\While{$b \leq \max\_bt$}
\State $W^{(t+1)}_{\text{cand}} = W^{(t)} \odot [(\frac{1}{\lambda} + \frac{1}{\theta} \max(|a| - \tau, 0))$ is the soft
\State \textbf{(Column normalization)} operator, and $w(\tilde{\eta})$, $v(\tilde{\eta})$ are standard
\State \hspace{0.5cm} Compute $D = \text{diag}(\frac{1}{\|w(\tilde{\eta})\|_2}, \dots, \frac{1}{\|w(\tilde{\eta})\|_2})$
\State \hspace{0.5cm} and set $W^{(t+1)}_{\text{cand}} = W^{(t+1)}_{\text{cand}} D$
\If{$\mathcal{L}(W^{(t+1)}_{\text{cand}}, H^{(t+1)}) \leq \mathcal{L}(W^{(t)}, H^{(t)})$}
\State $W^{(t+1)} = W^{(t+1)}_{\text{cand}}$
\State $H^{(t+1)} = H^{(t+1)}_{\text{cand}}$
\State $\beta^{(t)} = \beta^{(t)}_{\text{cand}}$ \subsubsection{Normalization of $W$}
\State $\text{flag\_accept} = \text{TRUE}$
\State break
\EndIf
\State $\theta = \theta * \rho$
\State $b = b + 1$
\EndWhile
\If{$\text{flag\_accept} = \text{FALSE}$}
\State $W^{(t+1)} = W^{(t)}$
\EndIf
\end{algorithmic}
\end{algorithm}

```

The column normalization preserves both \$WH\$ and \$W\beta\$, and therefore leaves the loss \$\mathcal{L}\$ invariant up to numerical errors. Subsection B.4.1. During the \$W\$ update from projected coordinate descent, the accepted \$W\$ update does not increase the \$\mathcal{L}\$.

Subsubsection{Update for \$\beta\$}

Recall that the overall loss function is

$$\mathcal{L}(W, H, \beta) = \frac{(1-\alpha)}{2np} \|X - WH\|^2 + \frac{\alpha}{2np} \sum_{i=1}^n \sum_{j=1}^n \beta_{ij} |w_{ij}|$$

Conditional on \$(W, H)\$, the loss in \$\beta\$ reduces to a convex elastic-net-penalized Cox problem:

$$\min_{\beta \in \mathbb{R}^k} \left\{ -\frac{2\alpha}{n} \sum_{i=1}^n \sum_{j=1}^k \beta_{ij} \ell_{ij}(W, \beta) + \lambda \sum_{i=1}^n \sum_{j=1}^k \beta_{ij}^2 \right\}$$

We solve this subproblem by cyclic coordinate descent following step size \$\gamma\$ be defined as in CITE: [Simon2011regularization]. Writing \$\ell(\beta) = \ell(W, \beta)\$, the update for coordinate \$i\$ in closed form

$$\hat{\beta}_i = \frac{S_i}{\frac{1}{\lambda} + \frac{1}{\alpha} \sum_{j=1}^n w_{ij}^2}$$

Then the update becomes

$$W^{(t)} \leftarrow W^{(t-1)} - \gamma \left(\frac{\partial \mathcal{L}}{\partial W} \right)_{t-1}$$

Finally, projected coordinate descent projects the W update in to the positive space

$$W^{(t)} = \max \left(W \odot \frac{\frac{(1-\alpha)}{np} X^T X + \frac{2\alpha}{n_{\text{event}}} \nabla_W \ell(\frac{(1-\alpha)}{np} X^T X + \frac{2\alpha}{n_{\text{event}}} \nabla_W \ell)}{\frac{(1-\alpha)}{np} X^T X + \frac{2\alpha}{n_{\text{event}}} \nabla_W \ell} \right)$$

Since $W \geq 0$ this is equivalent to

$$W^{(t)} = W \odot \max \left(\frac{\frac{(1-\alpha)}{np} X^T X + \frac{2\alpha}{n_{\text{event}}} \nabla_W \ell}{\frac{(1-\alpha)}{np} X^T X + \frac{2\alpha}{n_{\text{event}}} \nabla_W \ell}, 1 \right)$$

which matches the multiplicative form used in the software implementation of [simon2011regularization]

subproblem until convergence. Thus

Convergence proof

We show that, under mild regularity conditions, the block coordinate method (BCD) Algorithm-S\ref{alg:desurv} converges to a stationary point of the DeSurv loss function

$$\mathcal{L}(W, H, \beta)$$

For clarity, we first analyze the algorithm without the column normalization of W inside the loop, and then argue in Section~\ref{subsec:conv_normalization} that the normalization step preserves stationarity of limit points.

Throughout, let $\Theta = (W, H, \beta)$ and denote $\mathcal{L}(\Theta) = \mathcal{L}(W, H, \beta)$. The feasible set is

$$\Theta = \{(W, H, \beta) : W \in \mathbb{R}^{p \times k}, H \in \mathbb{R}^{k \times n}, \beta \in \mathbb{R}^k\}$$

Lemma 1 (Continuity and bounded level sets). Under Assumption 1, $\mathcal{L} : \Theta \rightarrow \mathbb{R}$ is continuous and its level sets are nonempty, closed, and bounded.

Assumption 1 (Regularity and parameter space). The NMF term $\|X - WH\|_F^2$ is a polynomial in the entries of W and H , and is continuously differentiable. The penalty $\|H\|_F^2$ is continuous and differentiable. The Cox partial log-likelihood is a smooth function of the linear coefficients β , and is continuously differentiable. The hyperparameters satisfy $\lambda_H > 0$, $\lambda_W > 0$, and $\lambda_\beta > 0$. The initial iterate $\Theta^{(0)} = (W^{(0)}, H^{(0)}, \beta^{(0)})$ is nonempty, closed, and bounded.

Assumption 2 (Block updates). The constraints $W, H \geq 0$ define closed convex cones, and β is unconstrained. Because $\lambda_H > 0$ and $\lambda_W > 0$, the terms $\frac{\lambda_H}{n} \|H\|_F^2$ and $\frac{\lambda_W}{n} \|W\|_F^2$ dominate the objective as $\|W\|_F, \|H\|_F \rightarrow \infty$, respectively. Thus $\mathcal{L}^{(t)} \mapsto \mathcal{L}^{(t+1)}$ is given by the multiplicative update rule $H \leftarrow H \odot \frac{X^T W}{X^T W + \lambda_H H}$ and $W \leftarrow W \odot \frac{X H}{X H + \lambda_W W}$. This rule preserves nonnegativity and yields a nonincreasing value of the reconstruction error $\|X - WH\|_F^2$ and $\mathcal{L}^{(t)}$; see e.g. [seung2001algorithms; pascualmontano2006nonsmooth].

Lemma 2 (Monotone descent and existence of limit points). Under Assumptions 1 and 2, Algorithm-S\ref{alg:desurv} (without backtracking as in Algorithm-S\ref{alg:backtrack}) generates a sequence $\{\Theta^{(t)}\}_{t \geq 0}$ such that

$$\mathcal{L}(W^{(t+1)}, H^{(t+1)}, \beta^{(t+1)}) \leq \mathcal{L}(\Theta^{(t+1)}) \leq \mathcal{L}(\Theta^{(t)})$$

(i) Conditional on (W, β) , the H -subproblem is convex in θ . Since θ^t converges to a finite limit θ^* , the sequence θ^t is bounded and therefore admits a limit point.

Proof.
By Assumption 2(i),

$$\mathcal{L}(W^{t+1}, H^{t+1}, \beta^t) \leq \mathcal{L}(W^t, H^t, \beta^t).$$
By Assumption 2(ii),

$$\mathcal{L}(W^{t+1}, H^{t+1}, \beta^t) \leq \mathcal{L}(W^t, H^{t+1}, \beta^t).$$
By Assumption 2(iii),

$$\mathcal{L}(W^{t+1}, H^{t+1}, \beta^{t+1}) \leq \mathcal{L}(W^{t+1}, H^{t+1}, \beta^t).$$
Combining these inequalities gives

$$\mathcal{L}(\theta^{t+1}) \leq \mathcal{L}(\theta^t).$$
Since $\mathcal{L}(\theta)$ is monotonically nonincreasing and θ is bounded below on Θ , the sequence θ^t converges to a finite limit θ^* .

Each iterate lies in \mathcal{S}_0 by construction, and \mathcal{S}_0 is bounded by Lemma 1. Therefore θ^t is bounded and has at least one limit point by the Bolzano-Weierstrass theorem. **Lemma 1 (Convergence to a stationary point).**
Under Assumptions 1 and 2, the sequence θ^t produced by Algorithm S (without W normalization) satisfies the following properties:
Let $\theta^* = (W^*, H^*, \beta^*)$ be any limit point of θ^t under Assumptions 1-2. Then:

- (i) H^* satisfies the KKT conditions for

$$\min_{H \geq 0} \mathcal{L}(W^*, H, \beta^*).$$
- (ii) W^* satisfies the KKT conditions for

$$\min_{W \geq 0} \mathcal{L}(W, H^*, \beta^*).$$
- (iii) β^* is the unique minimizer of

$$\min_{\beta \in \mathbb{R}^k} \mathcal{L}(W^*, H^*, \beta)$$
and satisfies the KKT conditions for the elastic net problem

$$\min_{\beta} \mathcal{L}(\theta) = f(W, H, \beta) + g(\beta) + I_{\{W \geq 0\}}(W) + I_{\{H \geq 0\}}(H),$$
where f is differentiable, $g(\beta) = \lambda \sum_i |\beta_i|$ is a separable convex (possibly nondifferentiable) penalty, and I_C is the indicator of a closed convex set C . This is the smooth+nonsmooth decomposition of $\mathcal{L}(\theta)$.

Proof.
Part (a) is Lemma 2. For part (b), Lemma 3 shows that any limit point is blockwise optimal: each block is optimal (in the KKT sense) for the subproblem with the other blocks fixed. The subproblem can be written as

$$\min_{\beta} \mathcal{L}(\theta) = f(W, H, \beta) + g(\beta)$$
which is differentiable, $g(\beta) = \lambda \sum_i |\beta_i|$ is a separable convex (possibly nondifferentiable) penalty, and I_C is the indicator of a closed convex set C . This is the smooth+nonsmooth decomposition of $\mathcal{L}(\theta)$.

form considered in block coordinate descent analyses such as [Weng 2006](#) and [Lange 2008](#) is generated to $((W, R^{-1})^T H)$ without change. In this setting, any point that is optimal with respect to a ℓ_1 block error, yielding multiple valid solutions individually (a block coordinatewise minimizer) is a stationary point of the overall problem. [Edonary 2004nmf](#), [Langeberg 2008](#) uniqueness; [Gillis 2014nmf](#). Although the full problem: equivalently, columns of (W) or rows of (H) resolves the trivial scaling and nonnegative mixing transformations persist, and empirical studies show that the sample-loading matrix (H) is more sensitive to its and local minima than the program matrix (W) [Brunet 2004metagenes](#); [Kim 2007sparse](#); [Gillis 2012accelerated](#)].

$$\begin{aligned} & \nabla_{\beta} f(W, H, \beta) \\ &= \nabla_{\beta} f(W, H, \beta) \\ &+ \partial_{\beta} I_{\{W \geq 0\}}(W) \\ &+ \partial_{\beta} I_{\{H \geq 0\}}(H), \end{aligned}$$

which are precisely the KKT conditions for $\min_{\theta \in \Theta} \mathcal{L}(\theta)$. Hence every limit point of Algorithm [S](#) [desurv](#) is a stationary point of the problem.

Effect of column normalization of W
 The above analysis omits the column-normalization step of Algorithm [S](#) [backtrack](#). We briefly argue that this step does not affect stationarity of limit points.

Let D be a diagonal matrix with strictly positive diagonal entries. Define the transformation

$$(W, H, \beta) \mapsto (W', H', \beta') = (W D^{-1}, D H, \beta)$$

preserves both the product WH and the linear predictor $\beta^T W$. If (W, H, β) is a limit point, then (W', H', β') is also a limit point. The column-normalization step of Algorithm [S](#) [backtrack](#) is a transformation, with D chosen from the column norms of W .

Let $\{\theta_t\}$ be the sequence generated by the algorithm, and let $\{\tilde{\theta}_t\}$ be the sequence generated by the algorithm with column normalization. For each t there exists a diagonal D_t with positive entries such that

$$(\tilde{\theta}_t)^T D_t^{-1} = \theta_t^T$$

and $\mathcal{L}(\tilde{\theta}_t) = \mathcal{L}(\theta_t)$. Thus limit points of $\{\tilde{\theta}_t\}$ are obtained from limit points of $\{\theta_t\}$ by such invertible diagonal scaling.

Since the KKT conditions are expressed in terms of the gradients with respect to WH and $W\beta$, and these quantities are invariant under the transformation, stationarity is preserved under the transformation. Theorem [1](#) implies that every limit point of the [normalized](#) algorithm is also a stationary point of \mathcal{L} (up to this scaling). This establishes convergence of the implementation used in practice.

Remark (Coxnet implementation)
 Our implementation updates the β block using a Coxnet-style coordinate descent that relies on an approximate Hessian.

Supervision on (W) versus (H)
 Classical nonnegative matrix factorization (NMF) is non-identifiable: for any invertible matrix (R) with nonnegative entries, the factorization

`\State` Number of initializations R
`\State` Observed data (X, y, Δ)
`\State` Hyperparameters k, α, λ, ξ , and λ_{c}
`\Ensure` The cross-validated c-index

`\State` Divide subjects into F folds.
`\For` $\{f = 1$ to $F\}$
`\State` Split data into training and validation $(X_{(-f)}, y_{(-f)}, \Delta_{(-f)})$ and $(X_{(f)}, y_{(f)}, \Delta_{(f)})$
`\For` $\{r = 1$ to $R\}$
`\State` set.seed(r)
`\State` Apply Algorithm~\ref{alg:desurv} with inputs $(X_{(-f)}, y_{(-f)}, \Delta_{(-f)})$, and (k, α, λ)
`\State` Obtain \hat{W} and $\hat{\beta}$ as output from Algorithm~\ref{alg:desurv}
`\State` Compute the estimated linear predictor:

$$\hat{\eta} = X_{(f)}^T \hat{W} \hat{\beta}$$

`\State` Compute c-index, denoted $\hat{c}_{(f)r}$, according to Equation~\ref{eq:c-index}
`\EndFor`
`\State` \textbf{end loop over r }
`\State` Compute average c-index across initializations:

$$\hat{c}_{(f)} = \frac{1}{R} \sum_{r=1}^R \hat{c}_{(f)r}$$

`\EndFor`
`\State` \textbf{end loop over f }
`\State` Compute final cross-validated c-index:

$$\hat{c} = \frac{1}{F} \sum_{f=1}^F \hat{c}_{(f)}$$

`\State` \Return Final estimate \hat{c}

`\end{algorithmic}`
`\end{algorithm}`

`\subsection{Bayesian Optimization}`
`\subsection{PDAC Datasets}`
`\subsection{Bladder Datasets}`
`\subsection{Consensus Clustering}`
`\subsection{Survival Analysis}`
`\pnasbreak`
`# Versioning {.unnumbered}`

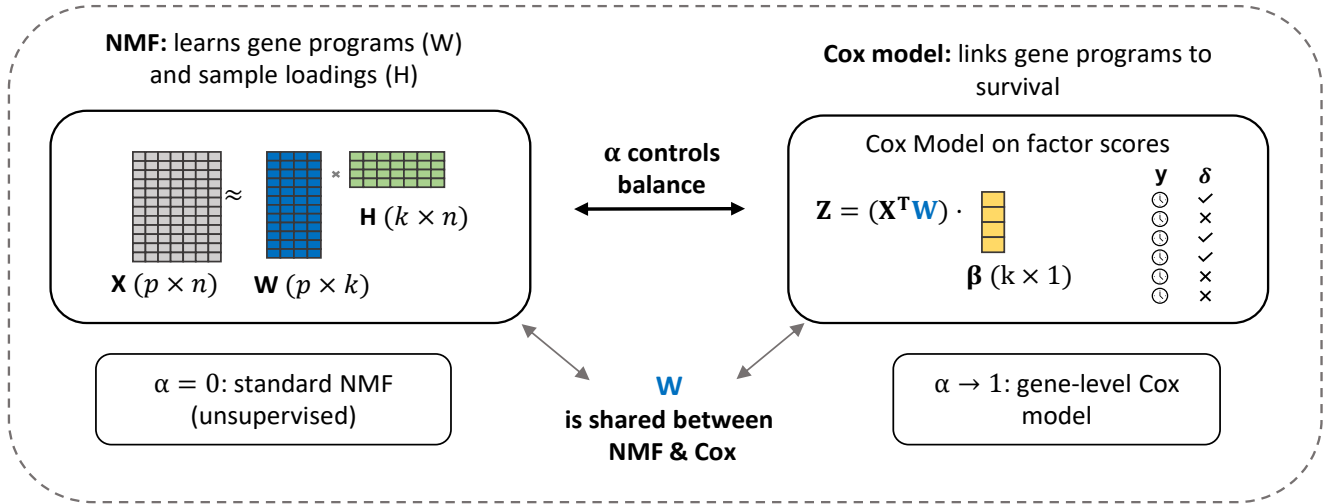
A. DeSurv package version: 1.0.1.
B. DeSurv git branch: 20260107bugfix.
C. DeSurv git commit: fea641a96e4743315a35b9b6addcc1a1ff53da9d.
D. Paper git branch: main.
E. Paper git commit: ad79f03a9361188583683b8f487de53d14c9f023.

- Pareja F, et al. (2016) Triple-negative breast cancer: The importance of molecular and histologic subtyping and recognition of low-grade variants. *NPJ breast cancer* 2(1):1–11.
- Dienstmann R, Salazar R, Tabernero J (2018) Molecular subtypes and the evolution of treatment decisions in metastatic colorectal cancer. *Am Soc Clin Oncol Educ Book* 38(38):231–8.
- Zhou X et al. (2021) Clinical impact of molecular subtyping of pancreatic cancer. *Frontiers in cell and developmental biology* 9:743908.
- Seiler R, et al. (2017) Impact of molecular subtypes in early-stage bladder cancer on clinical response and survival after neoadjuvant chemotherapy. *European urology* 72(4):544–554.
- Prat A, et al. (2015) Clinical implications of the intrinsic molecular subtypes of breast cancer. *The Breast* 24:S26–S35.
- Ou F, Michiels S, Shyr Y, Adjei AA, Oberg AL (2021) Biomarker discovery and validation: Statistical considerations. *Journal of Thoracic Oncology* 16(Suppl 15):S539–S547.
- Planey Catherine R, Gevaert O (2016) CoINcIDE: A framework for discovery of patient subtypes across multiple datasets. *Genome Medicine* 8(1):27.
- Prat A, Pineda E, Adamo B, et al. (2014) Molecular features and survival outcomes of the intrinsic subtypes in the international breast cancer study group trial 10-93. *Journal of the National Cancer Institute* 106(8):dju152.
- Ellrott K, et al. (2025) Classification of non-TCGA cancer samples to TCGA molecular subtypes using compact feature sets. *Cancer cell* 43(2):195–212.
- Tomczak K, Czerwińska P, Wiznerowicz M (2015) Review the cancer genome atlas (TCGA): An immeasurable source of knowledge. *Contemporary Oncology/Współczesna Onkologia* 2015(1):68–77.
- Zhang J, et al. (2019) The international cancer genome consortium data portal. *Nature biotechnology* 37(4):367–369.
- Nguyen H, Nguyen H, Tran D, Draghici S, Nguyen T (2024) Fourteen years of cellular deconvolution: Methodology, applications, technical evaluation and outstanding challenges. *Nucleic Acids Research* 52(9):4761–4783.
- Lee DD, Seung HS (1999) Learning the parts of objects by non-negative matrix factorization. *nature* 401(6755):788–791.

14. Bailey P, Chang DK, et al. (2016) [Genomic analyses identify molecular subtypes of pancreatic cancer](#). *Nature* 531(7592):47–52.
15. Moffitt RA, et al. (2015) Virtual microdissection identifies distinct tumor-and stroma-specific subtypes of pancreatic ductal adenocarcinoma. *Nature genetics* 47(10):1168–1178.
16. Peng XL, Moffitt RA, Torphy RJ, Volmar KE, Yeh JJ (2019) De novo compartment deconvolution and weight estimation of tumor samples using DECODER. *Nature communications* 10(1):4729.
17. Le Goff V, et al. (2025) SurvNMF: Non-negative matrix factorization supervised for survival data analysis. PhD thesis (Institut Pasteur Paris; CEA).
18. Huang Z, Salama P, Shao W, Zhang J, Huang K (2020) Low-rank reorganization via proportional hazards non-negative matrix factorization unveils survival associated gene clusters. *arXiv preprint arXiv:200803776*.

DRAFT

(A) The DeSurv model



(B) Data-driven model selection via Bayesian optimization

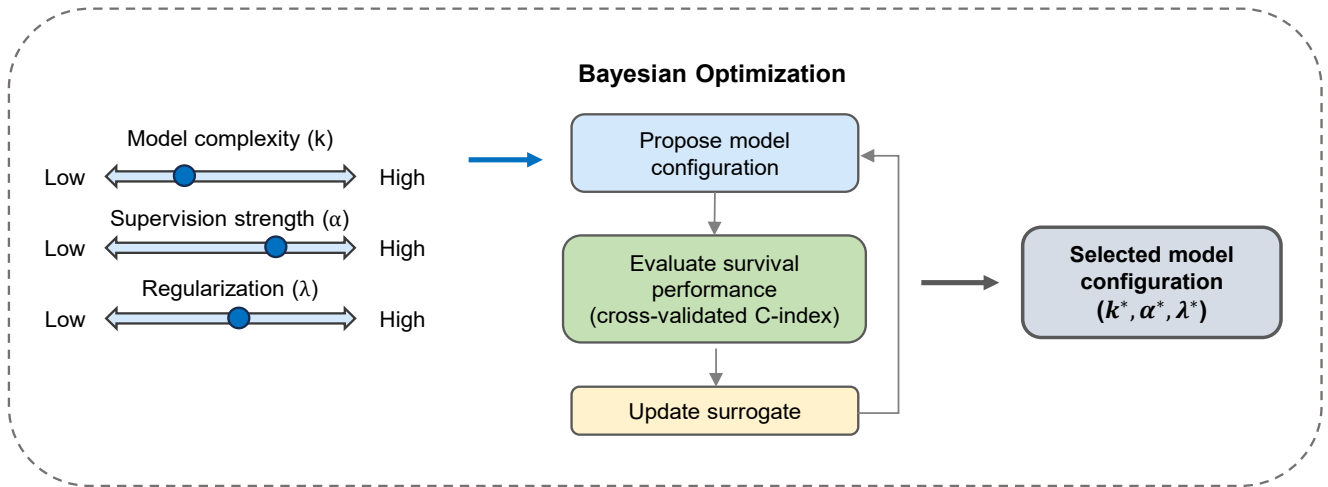


Fig. 1. Overview of the DeSurv framework and data-driven model selection. (A) DeSurv integrates nonnegative matrix factorization (NMF) with survival modeling to learn prognostic gene programs from a gene expression matrix X . NMF decomposes $X \approx WH$, where W represents gene programs and H sample loadings; the learned programs W are shared with a Cox proportional hazards model that links factor-derived scores $Z = X^T W$ to survival outcomes via regression coefficients β . A tuning parameter α controls the balance between unsupervised structure learning ($\alpha = 0$) and supervised survival association ($\alpha = 1$). (B) Model complexity (k), supervision strength (α), and regularization (λ) are selected via Bayesian optimization using cross-validated concordance index.